

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2024-525155  
(P2024-525155A)

(43)公表日 令和6年7月10日(2024.7.10)

(51)国際特許分類	F I	テーマコード(参考)
C 1 2 Q 1/02 (2006.01)	C 1 2 Q 1/02	4 B 0 6 3
C 1 2 Q 1/68 (2018.01)	C 1 2 Q 1/68	4 B 0 6 5
C 1 2 N 5/071(2010.01)	C 1 2 N 5/071	

審査請求 未請求 予備審査請求 未請求 (全139頁)

(21)出願番号	特願2023-577304(P2023-577304)	(71)出願人	520445473 フラッグシップ パイオニアリング イノベーションズ シックス, エルエルシー アメリカ合衆国マサチューセッツ州 0 2 1 4 2, ケンブリッジ, ケンブリッジ・パークウェイ 5 5, エイス・フロア
(86)(22)出願日	令和4年6月15日(2022.6.15)	(74)代理人	100079108 弁理士 稲葉 良幸
(85)翻訳文提出日	令和6年2月6日(2024.2.6)	(74)代理人	100109346 弁理士 大貫 敏史
(86)国際出願番号	PCT/US2022/033685	(74)代理人	100117189 弁理士 江口 昭彦
(87)国際公開番号	WO2022/266259	(74)代理人	100134120 弁理士 内藤 和彦
(87)国際公開日	令和4年12月22日(2022.12.22)	(72)発明者	ウォルフ, ファビアン アレクサンダー 最終頁に続く
(31)優先権主張番号	63/210,930		
(32)優先日	令和3年6月15日(2021.6.15)		
(33)優先権主張国・地域又は機関	米国(US)		
(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, 最終頁に続く		

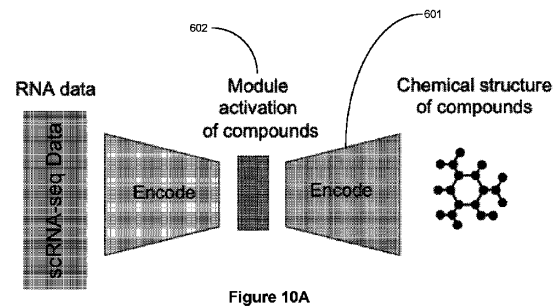
(54)【発明の名称】 フィンガープリント分析を使用して化合物を生理学的状態と関連付けるためのシステム及び方法

(57)【要約】

【課題】化合物を生理学的状態と関連付けるためのシステム及び方法が提供される。

【解決手段】化合物化学構造のフィンガープリントが得られ、1つ以上の計算された活性化スコアを出力するモデルに入力される。各活性化スコアは、モジュールのセットにおける細胞構成要素モジュールを表し、各モジュールは、細胞構成要素のサブセットを含み、モジュールのセットにおける第1のモジュールは、生理学的状態と関連付けられる。第1のモジュールについての活性化スコアが閾値基準を満たす場合、化合物は、生理学的状態と関連付けられると識別される。いくつかの態様において、各活性化スコアは、生理学的状態と関連付けられた摂動シグネチャを表し、化合物は、第1の摂動シグネチャについての活性化スコアが閾値基準を満たす場合に識別される。化合物を生理学的状態と関連付けるモデルを訓練するためのシステム及び方法も提供される。

【選択図】図10A



**【特許請求の範囲】****【請求項 1】**

試験化学化合物を目的の生理学的状態と関連付ける方法であって、前記方法が、

(A) 前記試験化学化合物の化学構造のフィンガープリントを得ることと、

(B) 細胞構成要素モジュールのセットにアクセスすることであって、

前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、

前記複数の細胞構成要素のそれぞれの独立したサブセットの各々についての対応する複数の細胞ベースのアッセイ存在量値が、前記生理学的状態と関連付けられた複数の異なる状態にわたって別々に相関し、

前記細胞構成要素モジュールのセットにおける第 1 の細胞構成要素モジュールが、前記目的の生理学的状態と関連付けられる、アクセスすることと、

(C) 前記化学構造の前記フィンガープリントを、100 以上のパラメータを含むモデルに入力することに応答して、前記モデルからの出力として、前記細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々についてのそれぞれの活性化スコアを取得することと、

(D) 前記第 1 の細胞構成要素モジュールについての前記活性化スコアが、第 1 の閾値基準を満たす場合、前記試験化学化合物を前記目的の生理学的状態と関連付けることとを含む、方法。

10

**【請求項 2】**

前記細胞ベースのアッセイ存在量値が、臓器の細胞のものである、請求項 1 に記載の方法。

20

**【請求項 3】**

前記臓器が、心臓、肝臓、肺、筋肉、脳、膵臓、脾臓、腎臓、小腸、子宮、又は膀胱である、請求項 2 に記載の方法。

**【請求項 4】**

前記細胞ベースのアッセイ存在量値が、組織の細胞のものである、請求項 1 に記載の方法。

**【請求項 5】**

前記組織が、骨、軟骨、関節、気管、脊髄、角膜、眼、皮膚、又は血管である、請求項 4 に記載の方法。

30

**【請求項 6】**

前記細胞ベースのアッセイ存在量値が、複数の幹細胞の細胞のものである、請求項 1 に記載の方法。

**【請求項 7】**

前記複数の幹細胞が、複数の胚性幹細胞、複数の成体幹細胞、又は複数の人工多能性幹細胞 (iPSC) である、請求項 6 に記載の方法。

**【請求項 8】**

前記細胞ベースのアッセイ存在量値が、複数の初代ヒト細胞の細胞のものである、請求項 1 に記載の方法。

40

**【請求項 9】**

前記複数の初代ヒト細胞が、複数の CD34+ 細胞、複数の CD34+ 造血幹、複数の前駆細胞 (HSPC)、複数の T 細胞、複数の間葉系幹細胞 (MSC)、複数の気道基底幹細胞、又は複数の人工多能性幹細胞である、請求項 8 に記載の方法。

**【請求項 10】**

前記細胞ベースのアッセイ存在量値が、臍帯血中、末梢血中、又は骨髓中の細胞のものである、請求項 1 に記載の方法。

**【請求項 11】**

前記細胞ベースのアッセイ存在量値が、固体組織中の細胞のものである、請求項 1 に記載の方法。

50

## 【請求項 1 2】

前記固体組織が、胎盤、肝臓、心臓、脳、腎臓、又は胃腸管である、請求項 1 1 に記載の方法。

## 【請求項 1 3】

前記細胞ベースのアッセイ存在量値が、複数の分化細胞のものである、請求項 1 に記載の方法。

## 【請求項 1 4】

前記複数の分化細胞が、複数の巨核球、複数の骨芽細胞、複数の軟骨細胞、複数の脂肪細胞、複数の肝細胞、複数の肝中皮細胞、複数の胆管上皮細胞、複数の肝星細胞、複数の肝類洞内皮細胞、複数のクッパー細胞、複数のピット細胞、複数の血管内皮細胞、複数の  
10  
膵管上皮細胞、複数の膵管細胞、複数の腺房中心細胞、複数の腺房細胞、複数のランゲルハンス島、複数の心筋細胞、複数の線維芽細胞、複数のケラチノサイト、複数の平滑筋細胞、複数の I 型肺胞上皮細胞、複数の II 型肺胞上皮細胞、複数のクララ細胞、複数の線毛上皮細胞、複数の基底細胞、複数の杯細胞、複数の神経内分泌細胞、複数のクルチッキー ( k u l t s c h i t z k y ) 細胞、複数の尿細管上皮細胞、複数の尿路上皮細胞、複数の円柱上皮細胞、複数の糸球体上皮細胞、複数の糸球体内皮細胞、複数の有足細胞、複数のメサングウム細胞、複数の神経細胞、複数の星状膠細胞、複数の小膠細胞、又は複数の乏突起膠細胞である、請求項 1 3 に記載の方法。

## 【請求項 1 5】

前記対応する複数の細胞ベースのアッセイ存在量値が、複数の細胞の単一細胞リボ核酸 ( R N A ) 配列決定 ( s c R N A - s e q ) データである、請求項 1 ~ 1 4 のいずれか一  
20  
項に記載の方法。

## 【請求項 1 6】

前記生理学的状態と関連付けられた前記複数の異なる状態が、細胞のアリコートが前記生理学的状態に影響を与えることが知られている化合物に曝露されている対照状態に加えて、前記生理学的状態に影響を与えることが知られている 1 つ以上の参照化合物に異なる細胞のアリコートを曝露することによって導出される、請求項 1 5 に記載の方法。

## 【請求項 1 7】

前記対応する複数の細胞ベースのアッセイ存在量値が、バルク R N A 配列に由来する、請求項 1 ~ 1 4 のいずれか一項に記載の方法。  
30

## 【請求項 1 8】

前記対応する複数の細胞ベースのアッセイ存在量値が、単一細胞 R N A 配列決定に由来する、請求項 1 ~ 1 4 のいずれか一項に記載の方法。

## 【請求項 1 9】

前記細胞構成要素モジュールのセットが、前記第 1 の細胞構成要素モジュールからなる、請求項 1 ~ 1 8 のいずれか一項に記載の方法。

## 【請求項 2 0】

前記細胞構成要素モジュールのセットが、複数の細胞構成要素モジュールを含み、前記モデルが、複数のコンポーネントモデルを含むアンサンブルモデルであり、前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、前記化学構造の前記フィン  
40  
ガープリントを前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々に入力することに対応して、前記細胞構成要素モジュールのセットにおける異なる細胞構成要素モジュールについての活性化スコアを提供する、請求項 1 ~ 1 8 のいずれか一項に記載の方法。

## 【請求項 2 1】

前記方法が、前記試験化学化合物の単純化された分子入力ラインエントリーシステム ( S M I L E S ) 文字列表現から前記フィンガープリントを計算することを更に含む、請求項 1 ~ 2 0 のいずれか一項に記載の方法。

## 【請求項 2 2】

前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、対応する二  
50

ューラルネットワークである、請求項 20 又は 21 に記載の方法。

【請求項 23】

前記対応するニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、22 に記載の方法。

【請求項 24】

前記複数のコンポーネントモデルにおけるコンポーネントモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルである、請求項 20 又は 21 に記載の方法。

10

【請求項 25】

前記対応するニューラルネットワークが、対応する完全に接続されたニューラルネットワーク及び対応するメッセージパッシングニューラルネットワークの組み合わせであり、前記対応する完全に接続されたニューラルネットワークの第 1 の出力及び前記対応するメッセージパッシングニューラルネットワークの第 2 の出力が、前記化学構造の前記フィンガープリントを前記対応する完全に接続されたニューラルネットワーク及び前記対応するメッセージパッシングニューラルネットワークに入力することに対応して、組み合わせられ、前記細胞構成要素モジュールのセットにおける前記対応する細胞構成要素モジュールについての 1 つ以上の計算された活性化スコアにおける活性化スコアを決定する、請求項 22 に記載の方法。

20

【請求項 26】

前記細胞構成要素モジュールのセットが、複数の細胞構成要素モジュールであり、前記第 1 の細胞構成要素モジュールを含む前記複数の細胞構成要素モジュールの第 1 のサブセットが、前記目的の生理学的状態と関連付けられ、前記複数の細胞構成要素モジュールの第 2 のサブセットが、前記目的の生理学的状態と関連付けられず、前記第 1 の細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが、前記第 1 の閾値基準を満たし、前記複数の細胞構成要素モジュールの前記第 2 のサブセットにおける細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが、前記第 1 の閾値基準以外の第 2 の閾値基準を満たす場合、前記試験化学化合物が、前記目的の生理学的状態と識別される、請求項 1 に記載の方法。

30

【請求項 27】

前記方法が、電子形式で 1 つ以上の第 1 のデータセットを得、前記 1 つ以上の第 1 のデータセットが、第 1 の複数の細胞におけるそれぞれの細胞の各々について、前記第 1 の複数の細胞が、20 個以上の細胞を含み、複数の注釈付きの細胞状態を集合的に表し、前記複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記複数の細胞構成要素が、10 個以上の細胞構成要素を含み、前記それぞれの細胞における前記それぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、それによって、複数のベクトルにアクセスするか、又はそれらを形成することによって、前記複数のベクトルにおけるそれぞれのベクトルの各々が、(i) 前記複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii) 対応する複数のエレメントを含み、前記対応する複数のエレメントにおけるそれぞれのエレメントの各々が、前記第 1 の複数の細胞における前記それぞれの細胞における前記それぞれの細胞構成要素の前記対応する存在量を表す対応するカウントを有する、複数のベクトルにアクセスするか、又はそれらを形成することと、

40

前記複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞

50

胞構成要素モジュールの各々を識別することであって、前記複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々が、前記複数の細胞構成要素のサブセットを含み、前記複数の細胞構成要素モジュールが、(i)前記複数の候補細胞構成要素モジュール及び(ii)前記複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、前記複数の細胞構成要素モジュールが、10を超える細胞構成要素モジュールを含む、識別することと、

電子形式で1つ以上の第2のデータセットを得、前記1つ以上の第2のデータセットが、

第2の複数の細胞におけるそれぞれの細胞の各々について、前記第2の複数の細胞が、20個以上の細胞を含み、前記目的の生理学的状態を通知する複数の共変量を集合的に表し、

10

前記複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、

前記それぞれの細胞における前記それぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、

それによって、(i)前記第2の複数の細胞及び(ii)前記複数の細胞構成要素又はその前記表現によって次元決定された細胞構成要素カウントデータ構造を得ることと、

前記複数の細胞構成要素又はその前記表現を共通次元として使用して前記細胞構成要素カウントデータ構造及び前記潜在表現を組み合わせることによって活性化データ構造を形成することであって、前記活性化データ構造が、前記複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、

20

前記第2の複数の細胞における細胞の各々について、それぞれの活性化重みを含む、形成することと、

前記複数の共変量におけるそれぞれの共変量の各々について、(i)前記共変量のフィンガープリントの候補細胞構成要素モデルへの入力時に、候補細胞構成要素モデルによって表される細胞構成要素モジュールの各々に対する計算された活性化と、(ii)前記候補細胞構成要素モデルによって表される細胞構成要素モジュールの各々に対する実際の活性化との間の差を使用して、前記候補細胞構成要素モデルを訓練することであって、前記訓練することが、前記差に応答して、前記候補細胞構成要素モデルと関連付けられた複数の共変量パラメータを調整する、訓練することと、を含む、プロセスによって前記第1の細胞構成要素モジュールを識別することを更に含む、請求項1~26のいずれか一項に記載の方法。

30

#### 【請求項28】

前記複数の共変量パラメータが、

前記複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、

それぞれの共変量の各々について、

前記それぞれの共変量が、前記第2の複数の細胞にわたって、前記それぞれの細胞構成要素モジュールと相関するかどうかを示す対応するパラメータを含み、前記方法が

前記候補細胞構成要素モデルを訓練する際に前記複数の共変量パラメータを使用して、前記複数の候補細胞構成要素モジュールにおける前記第1の細胞構成要素モジュールを識別することを更に含む、請求項27に記載の方法。

40

#### 【請求項29】

前記複数の注釈付きの細胞状態における注釈付きの細胞状態が、曝露条件下での化合物への前記第1の複数の細胞における細胞の曝露である、請求項27又は28に記載の方法。

#### 【請求項30】

前記曝露条件が、曝露期間、前記化合物の濃度、又は曝露期間及び前記化合物の濃度の組み合わせである、請求項29に記載の方法。

#### 【請求項31】

50

前記複数の細胞構成要素における細胞構成要素の各々が、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである、請求項 1 ~ 30 のいずれか一項に記載の方法。

【請求項 32】

前記複数の細胞構成要素における細胞構成要素の各々が、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせであり、

前記第 1 又は第 2 の複数の細胞における前記それぞれの細胞における前記それぞれの細胞構成要素の前記対応する存在量が、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動 (FRET) 測定によって決定される、請求項 27 ~ 30 のいずれか一項に記載の方法。

10

【請求項 33】

前記複数の細胞構成要素における細胞構成要素の各々が、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせであり、

前記第 1 又は第 2 の複数の細胞における前記それぞれの細胞における前記それぞれの細胞構成要素の前記対応する存在量が、単一細胞リボ核酸 (RNA) 配列決定 (scRNA-seq)、scTag-seq、配列決定を使用したトランスポザーゼ-アクセス可能なクロマチンのための単一細胞アッセイ (scATAC-seq)、CyTOF/SCOPE、EMSA/Abseq、miRNA-seq、CITE-seq、又はそれらの任意の組み合わせによって決定される、請求項 11 に記載の方法。

20

【請求項 34】

前記複数のベクトルを使用して、前記複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することが、前記複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、前記複数のベクトルに相関モデルを適用することを含む、請求項 1 ~ 30 又は 32 ~ 33 のいずれか一項に記載の方法。

【請求項 35】

前記相関モデルが、グラフクラスタリングを含む、請求項 34 に記載の方法。

【請求項 36】

前記グラフクラスタリングが、ピアソン相関ベースの距離メトリック上のライデン (Leiden) クラスタリングである、請求項 34 に記載の方法。

30

【請求項 37】

前記グラフクラスタリングが、ルーバン (Louvain) クラスタリングである、請求項 34 に記載の方法。

【請求項 38】

前記複数の細胞構成要素モジュールが、10 ~ 2000 個の細胞構成要素モジュールからなる、請求項 27 ~ 37 のいずれか一項に記載の方法。

【請求項 39】

前記複数の細胞構成要素が、100 ~ 8,000 個の細胞構成要素からなる、請求項 27 ~ 37 のいずれか一項に記載の方法。

40

【請求項 40】

前記複数の構成要素モジュールにおける候補細胞構成要素モジュールの各々が、200 ~ 300 個の細胞構成要素からなる、請求項 27 ~ 37 のいずれか一項に記載の方法。

【請求項 41】

前記目的の生理学的状態が、疾患である、請求項 1 ~ 40 のいずれか一項に記載の方法。

【請求項 42】

前記目的の生理学的状態が、疾患であり、前記第 1 の複数の細胞が、前記複数の注釈付きの細胞状態によって示されるように、前記疾患を代表する細胞、及び前記疾患を代表しない細胞を含む、請求項 27 に記載の方法。

50

## 【請求項 4 3】

前記複数の共変量が、細胞バッチ、細胞ドナー、細胞型、疾患状態、化学化合物への曝露、又はそれらの任意の組み合わせを含む、請求項 2 7 に記載の方法。

## 【請求項 4 4】

前記候補細胞構成要素モデルを前記訓練することが、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、前記複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、前記複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する、請求項 2 7 に記載の方法。

## 【請求項 4 5】

前記試験化学化合物が、2000ダルトン未満の分子量を有する有機化合物である、請求項 1 ~ 4 4 のいずれか一項に記載の方法。 10

## 【請求項 4 6】

前記試験化学化合物が、5つの基準のリピンスキーの法則の各々を満たす有機化合物である、請求項 4 5 に記載の方法。

## 【請求項 4 7】

前記試験化学化合物が、5つの基準の前記リピンスキーの法則のうちの少なくとも3つの基準を満たす有機化合物である、請求項 4 5 に記載の方法。

## 【請求項 4 8】

前記モデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項 1 ~ 1 9 のいずれか一項に記載の方法。 20

## 【請求項 4 9】

前記方法が、Daylight、BCI、ECFP4、ECFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNNS2S、又はGraphConvを使用して、前記試験化学化合物の化学構造から前記フィンガープリントを生成することを更に含む、請求項 1 ~ 4 8 のいずれか一項に記載の方法。

## 【請求項 5 0】

前記細胞構成要素モジュールのセットが、5つ以上の細胞構成要素モジュールを含む、請求項 1 ~ 1 8 又は 2 0 ~ 4 9 のいずれか一項に記載の方法。 30

## 【請求項 5 1】

前記細胞構成要素モジュールのセットが、10個以上の細胞構成要素モジュールを含む、請求項 1 ~ 1 8 又は 2 0 ~ 5 0 のいずれか一項に記載の方法。

## 【請求項 5 2】

前記細胞構成要素モジュールのセットが、100個以上の細胞構成要素モジュールを含む、請求項 1 ~ 1 8 又は 2 0 ~ 5 0 のいずれか一項に記載の方法。

## 【請求項 5 3】

前記それぞれの細胞構成要素モジュールにおける前記複数の細胞構成要素の前記独立したサブセットが、5つ以上の細胞構成要素を含む、請求項 1 ~ 5 2 のいずれか一項に記載の方法。 40

## 【請求項 5 4】

前記それぞれの細胞構成要素モジュールにおける前記複数の細胞構成要素の前記独立したサブセットが、前記目的の生理学的状態と関連付けられた分子経路における2 ~ 2 0 個の細胞構成要素からなる、請求項 1 ~ 5 2 のいずれか一項に記載の方法。

## 【請求項 5 5】

前記第1の閾値基準が、前記第1の細胞構成要素モジュールが閾値活性化スコアを有することが必要である、請求項 1 ~ 5 4 のいずれか一項に記載の方法。

## 【請求項 5 6】

1つ以上のプロセッサ及びメモリを含むコンピュータシステムであって、前記メモリが、試験化学化合物を目的の生理学的状態と関連付けるための方法を実施するための命令を 50

格納し、前記方法が、

(A) 前記試験化学化合物の化学構造のフィンガープリントを得ることと、

(B) 細胞構成要素モジュールのセットにアクセスすることであって、

前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、

前記複数の細胞構成要素のそれぞれの独立したサブセットの各々についての対応する複数の細胞ベースのアッセイ存在量値が、前記生理学的状態と関連付けられた複数の異なる状態にわたって別々に相関し、

前記細胞構成要素モジュールのセットにおける第1の細胞構成要素モジュールが、前記目的の生理学的状態と関連付けられる、アクセスすることと、

10

(C) 前記化学構造の前記フィンガープリントを、100以上のパラメータを含むモデルに入力することに対応して、前記モデルからの出力として、前記細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々についてのそれぞれの活性化スコアを取得することと、

(D) 前記第1の細胞構成要素モジュールについての前記活性化スコアが、第1の閾値基準を満たす場合、前記試験化学化合物を前記目的の生理学的状態と関連付けることと、を含む、コンピュータシステム。

【請求項57】

試験化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体であって、前記コンピュータが、1つ以上のプロセッサ及びメモリを含み、前記1つ以上のコンピュータプログラムが、

20

(A) 前記試験化学化合物の化学構造のフィンガープリントを得ることと、

(B) 細胞構成要素モジュールのセットにアクセスすることであって、

前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、

前記複数の細胞構成要素のそれぞれの独立したサブセットの各々についての対応する複数の細胞ベースのアッセイ存在量値が、前記生理学的状態と関連付けられた複数の異なる状態にわたって別々に相関し、

前記細胞構成要素モジュールのセットにおける第1の細胞構成要素モジュールが、前記目的の生理学的状態と関連付けられる、アクセスすることと、

30

(C) 前記化学構造の前記フィンガープリントを、100以上のパラメータを含むモデルに入力することに対応して、前記モデルからの出力として、前記細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々についてのそれぞれの活性化スコアを取得することと、

(D) 前記第1の細胞構成要素モジュールについての前記活性化スコアが、第1の閾値基準を満たす場合、前記試験化学化合物を前記目的の生理学的状態と関連付けることと、を含む、方法を実行するコンピュータによって実行可能な命令を集合的に符号化する、非一時的なコンピュータ可読媒体。

【請求項58】

40

試験化学化合物を目的の生理学的状態と関連付ける方法であって、前記方法が、

(A) 前記試験化学化合物の化学構造のフィンガープリントを得ることと、

(B) 摂動シグネチャのセットにアクセスすることであって、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、前記それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含む、前記それぞれの第1の細胞状態及び第2

50



の細胞状態のうち的一方が、非摂動細胞状態であり、前記それぞれの第1の細胞状態及び前記第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、アクセスすることと、

(C) 前記フィンガープリントをモデルに入力することであって、

前記モデルが、100以上のパラメータを含み、

前記モデルが、前記フィンガープリントの前記モデルへの前記入力にตอบสนองして1つ以上の計算された活性化スコアを出力し、

前記1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々が、前記摂動シグネチャのセットにおける対応する摂動シグネチャを表す、入力することと、

(D) 前記摂動シグネチャのセットにおける第1の摂動シグネチャについての前記それぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、前記化学化合物を前記目的の生理学的状態と関連付けることと、を含む、方法。

10

【請求項59】

前記方法が、前記試験化学化合物の単純化された分子入力ラインエントリーシステム(SMILES)文字列表現から前記フィンガープリントを計算することを更に含む、請求項58に記載の方法。

【請求項60】

前記モデルが、ニューラルネットワークを含む、請求項58又は59に記載の方法。

【請求項61】

前記ニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、請求項60に記載の方法。

20

【請求項62】

前記モデルが、複数のコンポーネントモデルを含むアンサンブルモデルであり、前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、前記化学構造の前記フィンガープリントを複数のコンポーネントモデルのセットにおけるコンポーネントモデルの各々に入力することに対応して、前記摂動シグネチャのセットにおける異なる摂動シグネチャについての活性化スコアを提供する、請求項58～61のいずれか一項に記載の方法。

30

【請求項63】

前記複数のコンポーネントモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項62に記載の方法。

【請求項64】

前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、対応するニューラルネットワークである、請求項62又は63に記載の方法。

【請求項65】

前記対応するニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、請求項64に記載の方法。

40

【請求項66】

前記複数のコンポーネントモデルにおけるコンポーネントモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルである、請求項63又は64に記載の方法。

【請求項67】

前記対応するニューラルネットワークが、完全に接続されたニューラルネットワーク及

50

びメッセージパッシングニューラルネットワークの組み合わせであり、

前記第1のニューラルネットワークの第1の出力及び前記第2のニューラルネットワークの第2の出力が、前記化学構造の前記フィンガープリントを前記完全に接続されたニューラルネットワーク及び前記メッセージパッシングニューラルネットワークに入力することに応答して、組み合わせられて、前記摂動シグネチャのセットにおける第1の摂動シグネチャについての前記1つ以上の計算された活性化スコアにおける活性化スコアを決定する、請求項65に記載の方法。

【請求項68】

前記摂動シグネチャのセットが、複数の摂動シグネチャであり、

前記第1の摂動シグネチャを含む、前記複数の摂動シグネチャの第1のサブセットが、前記目的の生理学的状態と関連付けられ、

前記複数の摂動シグネチャの第2のサブセットが、前記目的の生理学的状態と関連付けられておらず、

前記第1の摂動シグネチャについての前記それぞれの計算された活性化スコアが、前記第1の閾値基準を満たし、前記複数の摂動シグネチャの前記第2のサブセットにおける摂動シグネチャについての前記それぞれの計算された活性化スコアが、前記第1の閾値基準以外の第2の閾値基準を満たす場合、前記試験化学化合物が、前記目的の生理学的状態と識別される、請求項58に記載の方法。

【請求項69】

前記目的の生理学的状態が、疾患である、請求項58~68のいずれか一項に記載の方法。

【請求項70】

前記試験化学化合物が、2000ダルトン未満の分子量を有する有機化合物である、請求項58に記載の方法。

【請求項71】

前記試験化学化合物が、5つの基準のリピンスキーの法則の各々を満たす有機化合物である、請求項70に記載の方法。

【請求項72】

前記試験化学化合物が、5つの基準の前記リピンスキーの法則のうちの少なくとも3つの基準を満たす有機化合物である、請求項70に記載の方法。

【請求項73】

前記モデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項58に記載の方法。

【請求項74】

前記方法が、Daylight、BCI、ECFP4、EcFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNNS2S、又はGraphConvを使用して、前記試験化学化合物の化学構造から前記フィンガープリントを生成することを更に含む、請求項58~73のいずれか一項に記載の方法。

【請求項75】

前記摂動シグネチャのセットが、前記第1の摂動シグネチャからなる、請求項58~74のいずれか一項に記載の方法。

【請求項76】

前記摂動シグネチャのセットが、5つ以上の摂動シグネチャを含む、請求項58~74のいずれか一項に記載の方法。

【請求項77】

前記摂動シグネチャのセットが、10個以上の摂動シグネチャを含む、請求項58~74のいずれか一項に記載の方法。

【請求項78】

10

20

30

40

50

前記摂動シグネチャのセットが、100個以上の摂動シグネチャを含む、請求項58～74のいずれか一項に記載の方法。

【請求項79】

前記第1の閾値基準が、前記第1の摂動シグネチャが閾値活性化スコアを有することが必要である、請求項58～74のいずれか一項に記載の方法。

【請求項80】

1つ以上のプロセッサ及びメモリを含むコンピュータシステムであって、前記メモリが、試験化学化合物を目的の生理学的状態と関連付けるための方法を実施するための命令を格納し、前記方法が、

(A) 前記試験化学化合物の化学構造のフィンガープリントを得ることと、

10

(B) 摂動シグネチャのセットにアクセスすることであって、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、前記それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、前記それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、前記それぞれの第1の細胞状態及び前記第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、アクセスすることと、

20

(C) 前記フィンガープリントをモデルに入力することであって、

前記モデルが、100以上のパラメータを含み、

前記モデルが、前記フィンガープリントの前記モデルへの前記入力に応答して1つ以上の計算された活性化スコアを出力し、

前記1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々が、前記摂動シグネチャのセットにおける対応する摂動シグネチャを表す、入力することと、

(D) 前記摂動シグネチャのセットにおける第1の摂動シグネチャについての前記それぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、前記化学化合物を前記目的の生理学的状態と関連付けることと、を含む、コンピュータシステム。

30

【請求項81】

試験化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体であって、前記コンピュータが、1つ以上のプロセッサ及びメモリを含み、前記1つ以上のコンピュータプログラムが、

(A) 前記試験化学化合物の化学構造のフィンガープリントを得ることと、

(B) 摂動シグネチャのセットにアクセスすることであって、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、

40

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、前記それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、前記それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、前記それぞれの第1の細胞状態及び前記第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、アクセスすることと、

(C) 前記フィンガープリントをモデルに入力することであって、

前記モデルが、100以上のパラメータを含み、

50

前記モデルが、前記フィンガープリントの前記モデルへの前記入力に応答して1つ以上の計算された活性化スコアを出力し、

前記1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々が、前記摂動シグネチャのセットにおける対応する摂動シグネチャを表す、入力することと、

(D) 前記摂動シグネチャのセットにおける第1の摂動シグネチャについての前記それぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、前記化学化合物を前記目的の生理学的状態と関連付けることと、を含む、方法を実行するコンピュータによって実行可能な命令を集的に符号化する、非一時的なコンピュータ可読媒体。

【請求項82】

化学化合物を目的の生理学的状態と関連付ける方法であって、前記方法が、メモリ及び1つ以上のプロセッサを含むコンピュータシステムにおいて、

(A) 複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることと、

(B) 前記複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む、得ることと、

(C)

前記複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、

前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、

(i) 前記それぞれの化合物の前記化学構造の前記フィンガープリントを訓練されていないモデルに入力したときの前記それぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、(ii) 前記細胞構成要素モジュールのセットにおける前記それぞれの化合物についての前記それぞれの細胞構成要素モジュールの前記それぞれの数値的活性化スコアとの間のそれぞれの差を使用して前記訓練されていないモデルを訓練することであって、前記訓練すること(C)が、前記差に応答して前記訓練されていないモデルと関連付けられた複数のパラメータを調整し、前記複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を前記目的の生理学的状態と関連付ける訓練されたモデルを得る、訓練することと、を含む、方法。

【請求項83】

前記細胞構成要素モジュールのセットが、単一の細胞構成要素モジュールからなる、請求項82に記載の方法。

【請求項84】

前記細胞構成要素モジュールのセットが、複数の細胞構成要素モジュールを含む、請求項82に記載の方法。

【請求項85】

前記細胞構成要素モジュールのセットが、200~500個の細胞構成要素モジュールからなる、請求項82に記載の方法。

【請求項86】

前記複数の化合物が、 $10 \sim 1 \times 10^6$ 個の化合物からなる、請求項82に記載の方法。

【請求項87】

前記複数の化合物が、 $100 \sim 100,000$ 個の化合物からなる、請求項82に記載の方法。

【請求項88】

10

20

30

40

50

前記複数の化合物が、1000～100,000個の化合物からなる、請求項82に記載の方法。

【請求項89】

前記訓練すること(C)が、回帰アルゴリズムに従って、前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々についてのそれぞれの化合物の各々と関連付けられた差の各々に応答して、前記訓練されていないモデルと関連付けられた前記複数のパラメータを調整する、請求項82～88のいずれか一項に記載の方法。

【請求項90】

前記回帰アルゴリズムが、前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々についてのそれぞれの化合物の各々と関連付けられた差の各々の最小二乗誤差を最適化する、請求項89に記載の方法。

10

【請求項91】

前記訓練されたモデルが、ニューラルネットワークを含む、請求項82～90のいずれか一項に記載の方法。

【請求項92】

前記ニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、請求項91に記載の方法。

【請求項93】

前記訓練されたモデルが、複数のコンポーネントモデルのアンサンブルモデルであり、前記複数のコンポーネントモデルにおけるそれぞれのコンポーネントモデルの各々が、前記複数の細胞構成要素モジュールにおける異なる細胞構成要素モジュールについて計算された活性化スコアを出力する、請求項82～90のいずれか一項に記載の方法。

20

【請求項94】

前記複数のコンポーネントモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項93に記載の方法。

【請求項95】

前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、対応するニューラルネットワークである、請求項93に記載の方法。

30

【請求項96】

前記対応するニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、請求項95に記載の方法。

【請求項97】

前記細胞構成要素モジュールのセットが、複数の細胞構成要素モジュールであり、前記複数の細胞構成要素モジュールの第1のサブセットが、前記目的の生理学的状態と関連付けられ、

前記複数の細胞構成要素モジュールの第2のサブセットが、前記目的の生理学的状態と関連付けられていない、請求項82～96のいずれか一項に記載の方法。

40

【請求項98】

前記方法が、  
電子形式で1つ以上の第1のデータセットを得、前記1つ以上の第1のデータセットが、

第1の複数の細胞におけるそれぞれの細胞の各々について、前記第1の複数の細胞が、20個以上の細胞を含み、複数の注釈付きの細胞状態を集合的に表し、

前記複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記複数の細胞構成要素が、10個以上の細胞構成要素を含み、

前記それぞれの細胞における前記それぞれの細胞構成要素の対応する存在量を含

50

むか、又は集合的に含み、

それによって、複数のベクトルにアクセスするか、又はそれらを形成することによって、前記複数のベクトルにおけるそれぞれのベクトルの各々が、(i)前記複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii)対応する複数のエレメントを含み、前記対応する複数のエレメントにおけるそれぞれのエレメントの各々が、前記第1の複数の細胞における前記それぞれの細胞における前記それぞれの細胞構成要素の前記対応する存在量を表す対応するカウントを有する、複数のベクトルにアクセスするか、又はそれらを形成することと、

前記複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することによって、前記複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々が、前記複数の細胞構成要素のサブセットを含み、前記複数の細胞構成要素モジュールが、(i)前記複数の候補細胞構成要素モジュール及び(ii)前記複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、前記複数の細胞構成要素モジュールが、10を超える細胞構成要素モジュールを含む、識別することと、

電子形式で1つ以上の第2のデータセットを得、前記1つ以上の第2のデータセットが、

第2の複数の細胞におけるそれぞれの細胞の各々について、前記第2の複数の細胞が、20個以上の細胞を含み、前記目的の生理学的状態を通知する複数の共変量を集合的に表し、

前記複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、

前記それぞれの細胞における前記それぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、

それによって、(i)前記第2の複数の細胞及び(ii)前記複数の細胞構成要素又はその前記表現によって次元決定された細胞構成要素カウントデータ構造を得ることと、

前記複数の細胞構成要素又はその前記表現を共通次元として使用して前記細胞構成要素カウントデータ構造及び前記潜在表現を組み合わせることによって活性化データ構造を形成することによって、前記活性化データ構造が、前記複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、

前記第2の複数の細胞における細胞の各々について、それぞれの活性化重みを含む、形成することと、

(i)前記活性化データ構造を候補モデルに入力したときに、前記活性化データ構造内に表される細胞構成要素モジュールの各々における前記複数の共変量における各共変量の不在又は存在の予測と、(ii)細胞構成要素モジュールの各々における各共変量の実際の不在又は存在との間の差を使用して、候補細胞構成要素モデルを訓練することによって、前記訓練することが、前記差に応答して、前記候補細胞構成要素モデルと関連付けられた複数の共変量パラメータを調整する、訓練することと、を含む、プロセスによって前記複数の細胞構成要素モジュールにおける細胞構成要素モジュールを識別することを更に含む、請求項82~97のいずれか一項に記載の方法。

【請求項99】

前記複数の共変量パラメータが、

前記複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、

それぞれの共変量の各々について、

前記それぞれの共変量が、前記第2の複数の細胞にわたって、前記それぞれの細胞構成要素モジュールと相関するかどうかを示す対応するパラメータを含み、

前記候補細胞構成要素モデルを訓練する際に前記複数の共変量パラメータを使用して、前記複数の候補細胞構成要素モジュールにおける前記細胞構成要素モジュールを識別する、請求項98に記載の方法。

【請求項100】

10

20

30

40

50

前記複数の注釈付きの細胞状態における注釈付きの細胞状態が、曝露条件下での化合物への前記第1の複数の細胞における細胞の曝露である、請求項99に記載の方法。

【請求項101】

前記曝露条件が、曝露期間、前記化合物の濃度、又は曝露期間及び前記化合物の濃度の組み合わせである、請求項99に記載の方法。

【請求項102】

前記複数の細胞構成要素における細胞構成要素の各々が、特定の遺伝子、遺伝子に関連する特定のmRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである、請求項82～101のいずれか一項に記載の方法。

【請求項103】

前記第1又は第2の複数の細胞における前記それぞれの細胞における前記それぞれの細胞構成要素の前記対応する存在量が、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動(FRET)測定によって決定される、請求項98に記載の方法。

【請求項104】

前記第1又は第2の複数の細胞における前記それぞれの細胞における前記それぞれの細胞構成要素の前記対応する存在量が、単一細胞リボ核酸(RNA)配列決定(scRNA-seq)、scTag-seq、配列決定を使用したトランスポザーゼ-アクセス可能なクロマチンのための単一細胞アッセイ(scATAC-seq)、CyTOF/SCOPE、E-MS/Abseq、miRNA-seq、CITE-seq、又はそれらの任意の組み合わせによって決定される、請求項98に記載の方法。

【請求項105】

前記複数のベクトルを使用して、前記複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することが、前記複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、前記複数のベクトルに相関モデルを適用することを含む、請求項98に記載の方法。

【請求項106】

前記相関モデルが、グラフクラスタリングを含む、請求項105に記載の方法。

【請求項107】

前記グラフクラスタリング方法が、ピアソン相関ベースの距離メトリック上のライデン(Leiden)クラスタリングであるか、又はルーバン(Louvain)クラスタリングである、請求項106に記載の方法。

【請求項108】

前記複数の細胞構成要素が、100～8,000個の細胞構成要素からなる、請求項82～107のいずれか一項に記載の方法。

【請求項109】

前記複数の構成要素モジュールにおける候補細胞構成要素モジュールの各々が、200～300個の細胞構成要素からなる、請求項98に記載の方法。

【請求項110】

前記目的の生理学的状態が、疾患である、請求項82～109のいずれか一項に記載の方法。

【請求項111】

前記生理学的状態が、疾患であり、前記第1の複数の細胞が、前記複数の注釈付きの細胞状態によって示されるように、前記疾患を代表する細胞、及び前記疾患を代表しない細胞を含む、請求項98のいずれか一項に記載の方法。

【請求項112】

前記複数の共変量が、細胞バッチ、細胞ドナー、細胞型、疾患状態、又は化学化合物への曝露を含む、請求項98に記載の方法。

【請求項113】

前記候補細胞構成要素モデルを前記訓練することが、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、前記複数の共変量における共変量の各々が

10

20

30

40

50

、複数のコスト関数におけるコスト関数に対応し、前記複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する、請求項 98 に記載の方法。

【請求項 114】

前記複数の化学化合物における化学化合物の各々が、2000 ダルトン未満の分子量を有する有機化合物である、請求項 82 ~ 113 のいずれか一項に記載の方法。

【請求項 115】

前記複数の化学化合物における化学化合物の各々が、5つの基準のリピンスキーの法則の各々を満たす、請求項 82 ~ 113 のいずれか一項に記載の方法。

【請求項 116】

前記複数の化学化合物における化学化合物の各々が、5つの基準の前記リピンスキーの法則のうち少なくとも3つの基準を満たす、請求項 82 ~ 113 のいずれか一項に記載の方法。

【請求項 117】

前記訓練されたモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項 82 ~ 116 のいずれか一項に記載の方法。

【請求項 118】

前記方法が、Daylight、BCI、ECFP4、EcFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNNS2S、又はGraphConvを使用して、前記対応する化学構造からそれぞれのフィンガープリントの各々を生成することを更に含む、請求項 82 ~ 117 のいずれか一項に記載の方法。

【請求項 119】

前記細胞構成要素モジュールのセットが、5つ以上の細胞構成要素モジュールを含む、請求項 82 に記載の方法。

【請求項 120】

前記細胞構成要素モジュールのセットが、10個以上の細胞構成要素モジュールを含む、請求項 82 に記載の方法。

【請求項 121】

前記細胞構成要素モジュールのセットが、100個以上の細胞構成要素モジュールを含む、請求項 82 に記載の方法。

【請求項 122】

1つ以上のプロセッサ及びメモリを含むコンピュータシステムであって、前記メモリが、化学化合物を目的の生理学的状態と関連付けるための方法を実施するための命令を格納し、前記方法が、

(A) 複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることと、

(B) 前記複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む、得ることと、

(C)

前記複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、

前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、

(i) 前記それぞれの化合物の前記化学構造の前記フィンガープリントを訓練



されていないモデルに入力したときの前記それぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、( i i ) 前記細胞構成要素モジュールのセットにおける前記それぞれの化合物についての前記それぞれの細胞構成要素モジュールの前記それぞれの数値的活性化スコアとの間のそれぞれの差を使用して前記訓練されていないモデルを訓練することであって、前記訓練すること ( C ) が、前記差に応答して前記訓練されていないモデルと関連付けられた複数のパラメータを調整し、前記複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を前記目的の生理学的状態と関連付ける訓練されたモデルを得ることと、を含む、コンピュータシステム。

【請求項123】

化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体であって、前記コンピュータが、1つ以上のプロセッサ及びメモリを含み、前記1つ以上のコンピュータプログラムが、

10

( A ) 複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることと、

( B ) 前記複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む、得ることと、

20

( C )

前記複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、

前記細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、

( i ) 前記それぞれの化合物の前記化学構造の前記フィンガープリントを訓練されていないモデルに入力したときの前記それぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、( i i ) 前記細胞構成要素モジュールのセットにおける前記それぞれの化合物についての前記それぞれの細胞構成要素モジュールの前記それぞれの数値的活性化スコアとの間のそれぞれの差を使用して前記訓練されていないモデルを訓練することであって、前記訓練すること ( C ) が、前記差に応答して前記訓練されていないモデルと関連付けられた複数のパラメータを調整し、前記複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を前記目的の生理学的状態と関連付ける訓練されたモデルを得ることと、を含む、方法を実行するコンピュータによって実行可能な命令を集合的に符号化する、非一時的なコンピュータ可読媒体。

30

【請求項124】

化学化合物を目的の生理学的状態と関連付ける方法であって、前記方法が、

メモリ及び1つ以上のプロセッサを含むコンピュータシステムにおいて、

( A ) 複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることと、

40

( B ) 前記複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、前記それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、前記それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、前記それぞれの

50

第 1 の細胞状態及び前記第 2 の細胞状態のうちの方が、前記対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、得ることと、

( C )

前記複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について、

( i ) 前記それぞれの化合物の前記化学構造の前記フィンガープリントを訓練されていないモデルに入力したときの前記それぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、( i i ) 前記摂動シグネチャのセットにおける前記対応する化合物についての前記それぞれの摂動シグネチャの前記それぞれの数値的活性化スコアとの間のそれぞれの差を使用して前記訓練されていないモデルを訓練することであって、前記訓練すること( C )が、前記差に応答して、前記訓練されていないモデルと関連付けられた複数のパラメータを調整し、前記複数のパラメータが、100 以上のパラメータを含み、それによって、化学化合物を前記目的の生理学的状態と関連付ける訓練されたモデルを得る、訓練することと、を含む、方法。

10

【請求項 1 2 5】

前記摂動シグネチャのセットが、単一の摂動シグネチャからなる、請求項 1 2 4 に記載の方法。

【請求項 1 2 6】

前記摂動シグネチャのセットが、200 ~ 500 個の摂動シグネチャからなる、請求項 1 2 4 に記載の方法。

20

【請求項 1 2 7】

前記複数の化合物が、 $10 \sim 1 \times 10^6$  個の化合物からなる、請求項 1 2 4 ~ 1 2 6 のいずれか一項に記載の方法。

【請求項 1 2 8】

前記複数の化合物が、100 ~ 100,000 個の化合物からなる、請求項 1 2 4 ~ 1 2 6 のいずれか一項に記載の方法。

【請求項 1 2 9】

前記複数の化合物が、1000 ~ 100,000 個の化合物からなる、請求項 1 2 4 ~ 1 2 6 のいずれか一項に記載の方法。

30

【請求項 1 3 0】

前記訓練すること( C )が、回帰アルゴリズムに従って、前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々についての対応する化合物の各々と関連付けられた差の各々に応答して、前記訓練されていないモデルと関連付けられた前記複数のパラメータを調整する、請求項 1 2 4 ~ 1 2 9 のいずれか一項に記載の方法。

【請求項 1 3 1】

前記回帰アルゴリズムが、前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々についての対応する化合物の各々と関連付けられた差の各々の最小二乗誤差を最適化する、請求項 1 3 0 に記載の方法。

40

【請求項 1 3 2】

前記訓練されたモデルが、ニューラルネットワークを含む、請求項 1 2 4 ~ 1 3 1 のいずれか一項に記載の方法。

【請求項 1 3 3】

前記ニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、請求項 1 3 2 に記載の方法。

【請求項 1 3 4】

前記訓練されたモデルが、複数のコンポーネントモデルのアンサンブルモデルであり、前記複数のコンポーネントモデルにおけるそれぞれのコンポーネントモデルの各々が、そ

50

れぞれの化学構造のフィンガープリントを複数のコンポーネントモデルのセットにおけるコンポーネントモデルの各々に入力することに対応して、前記複数の摂動シグネチャのセットにおける異なる摂動シグネチャのセットについて計算された活性化スコアを出力する、請求項 1 2 4 に記載の方法。

【請求項 1 3 5】

前記複数のコンポーネントモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項 1 3 4 に記載の方法。

【請求項 1 3 6】

10

前記複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、対応するニューラルネットワークである、請求項 1 3 4 に記載の方法。

【請求項 1 3 7】

前記対応するニューラルネットワークが、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである、請求項 1 3 6 に記載の方法。

【請求項 1 3 8】

前記摂動シグネチャのセットが、複数の摂動シグネチャを含み、

前記複数の摂動シグネチャの第 1 のサブセットが、前記目的の生理学的状態と関連付けられ、

20

前記複数の摂動シグネチャの第 2 のサブセットが、前記目的の生理学的状態と関連付けられていない、請求項 1 2 4 ~ 1 3 7 のいずれか一項に記載の方法。

【請求項 1 3 9】

前記目的の生理学的状態が、疾患である、請求項 1 2 4 ~ 1 3 8 のいずれか一項に記載の方法。

【請求項 1 4 0】

前記複数の化学化合物における化学化合物の各々が、2000 ダルトン未満の分子量を有する有機化合物である、請求項 1 2 4 ~ 1 3 9 のいずれか一項に記載の方法。

【請求項 1 4 1】

前記複数の化学化合物における化学化合物の各々が、5 つの基準のリピンスキーの法則の各々を満たす、請求項 1 2 4 ~ 1 4 0 のいずれか一項に記載の方法。

30

【請求項 1 4 2】

前記複数の化学化合物における化学化合物の各々が、5 つの基準の前記リピンスキーの法則のうち少なくとも 3 つの基準を満たす、請求項 1 2 4 ~ 1 4 0 のいずれか一項に記載の方法。

【請求項 1 4 3】

前記訓練されたモデルが、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシン、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む、請求項 1 2 4 に記載の方法。

40

【請求項 1 4 4】

前記方法が、Daylight、BCI、ECFP4、EcFC、MDL、APFP、TTFP、UNITY 2D フィンガープリント、RNNS2S、又は GraphConv を使用して、前記対応する化学構造からそれぞれのフィンガープリントの各々を生成することを更に含む、請求項 1 2 4 ~ 1 4 3 のいずれか一項に記載の方法。

【請求項 1 4 5】

前記摂動シグネチャのセットが、5 つ以上の摂動シグネチャを含む、請求項 1 2 4 に記載の方法。

【請求項 1 4 6】

前記摂動シグネチャのセットが、10 個以上の摂動シグネチャを含む、請求項 1 2 4 に

50

記載の方法。

【請求項 1 4 7】

前記摂動シグネチャのセットが、100個以上の摂動シグネチャを含む、請求項 1 2 4 に記載の方法。

【請求項 1 4 8】

前記方法が、

変化していない細胞状態と変化した細胞状態との間の差次的細胞構成要素存在量の尺度を表す単一細胞遷移シグネチャに電子形式でアクセスすることであって、

前記変化した細胞状態が、前記変化していない細胞状態から前記変化した細胞状態への前記細胞遷移を通して発生し、

( i ) 前記変化していない細胞状態、( i i ) 前記変化した細胞状態、及び( i i i ) 前記変化していない細胞状態から前記変化した細胞状態への前記遷移のうちの少なくとも1つが、前記目的の生理学的状態と関連付けられ、

前記単一細胞遷移シグネチャが、参照の複数の細胞構成要素の識別と、前記複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、前記変化していない細胞状態と前記変化した細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する第1の有意性スコアとを含む、アクセスすることと、

前記単一細胞遷移シグネチャと前記それぞれの摂動シグネチャとを比較し、それによって前記それぞれの摂動シグネチャの前記それぞれの数値的活性化スコアを決定することと、を含む、手順によって前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアを得ることを更に含む、請求項 1 2 4 に記載の方法。

【請求項 1 4 9】

前記単一細胞遷移シグネチャと前記摂動シグネチャとを前記比較して、前記それぞれの摂動シグネチャの前記それぞれの数値的活性化スコアを決定することが、前記単一細胞遷移シグネチャの前記参照の複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、

前記それぞれの摂動シグネチャにおける前記対応する細胞構成要素の前記対応する有意性スコアに対する前記それぞれの細胞構成要素の前記第1の有意性スコアを比較することを含む、請求項 1 4 8 に記載の方法。

【請求項 1 5 0】

前記それぞれの摂動シグネチャの前記活性化スコアが、前記摂動シグネチャのセットにおける他の摂動シグネチャと比較して、前記単一細胞遷移シグネチャに対する前記それぞれの摂動シグネチャの関連性の相対的なランキングである、請求項 1 4 8 又は 1 4 9 に記載の方法。

【請求項 1 5 1】

前記相対的なランキングが、ウィルコクソンの順位和検定、t検定、ロジスティック回帰、又は一般化線形モデルによって決定される、請求項 1 5 0 に記載の方法。

【請求項 1 5 2】

前記単一細胞遷移シグネチャの前記変化していない細胞状態が、前記それぞれの摂動シグネチャの前記第1の細胞状態又は前記第2の細胞状態と同じである、請求項 1 4 8 ~ 1 5 1 のいずれか一項に記載の方法。

【請求項 1 5 3】

前記単一細胞遷移シグネチャの前記変化していない細胞状態が、前記それぞれの摂動シグネチャの前記第1の細胞状態及び前記第2の細胞状態の両方とは異なる、請求項 1 4 8 ~ 1 5 1 のいずれか一項に記載の方法。

【請求項 1 5 4】

前記方法が、

前記単一細胞遷移シグネチャの前記参照の複数の細胞構成要素、及び前記それぞれの摂動シグネチャの前記それぞれの複数の細胞構成要素を剪定して、転写因子と比較するこ

10

20

30

40

50

とを制限することを更に含む、請求項 1 4 8 ~ 1 5 3 のいずれか一項に記載の方法。

【請求項 1 5 5】

前記複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの前記摂動細胞状態が、前記複数の化合物における化合物に曝露されていない対照細胞によって表される、請求項 1 2 4 ~ 1 5 4 のいずれか一項に記載の方法。

【請求項 1 5 6】

前記複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの前記摂動細胞状態が、前記それぞれの摂動シグネチャと関連付けられた前記化合物以外の前記複数の化学化合物における化学化合物に曝露されている無関係の摂動細胞にわたる平均によって表される、請求項 1 2 4 ~ 1 5 4 のいずれか一項に記載の方法。

10

【請求項 1 5 7】

1 つ以上のプロセッサ及びメモリを含むコンピュータシステムであって、前記メモリが、化学化合物を目的の生理学的状態と関連付けるための命令を格納し、前記方法が、

( A ) 複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることと、

( B ) 前記複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、前記それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、それぞれの第 1 の細胞状態とそれぞれの第 2 の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、前記それぞれの第 1 の細胞状態及び第 2 の細胞状態のうち一方が、非摂動細胞状態であり、前記それぞれの第 1 の細胞状態及び前記第 2 の細胞状態のうち他方が、前記対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、得ることと、

20

( C )

前記複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について

30

( i ) 前記それぞれの化合物の前記化学構造の前記フィンガープリントを訓練されていないモデルに入力したときの前記それぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、( i i ) 前記摂動シグネチャのセットにおける前記対応する化合物についての前記それぞれの摂動シグネチャの前記それぞれの数値的活性化スコアとの間のそれぞれの差を使用して前記訓練されていないモデルを訓練することであって、前記訓練すること( C ) が、前記差に応答して、前記訓練されていないモデルと関連付けられた複数のパラメータを調整し、前記複数のパラメータが、1 0 0 以上のパラメータを含み、それによって、化学化合物を前記目的の生理学的状態と関連付ける訓練されたモデルを得る、訓練することと、を含む、コンピュータシステム。

40

【請求項 1 5 8】

化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な 1 つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体であって、前記コンピュータが、1 つ以上のプロセッサ及びメモリを含み、前記 1 つ以上のコンピュータプログラムが、

( A ) 複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることと、

( B ) 前記複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子

50

形式で得ることであって、前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、前記それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、前記それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、前記それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、前記それぞれの第1の細胞状態及び前記第2の細胞状態のうち他方が、前記対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、得ることと、

(C)

前記複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、

前記摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について

(i) 前記それぞれの化合物の前記化学構造の前記フィンガープリントを訓練されていないモデルに入力したときの前記それぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、(ii) 前記摂動シグネチャのセットにおける前記対応する化合物についての前記それぞれの摂動シグネチャの前記それぞれの数値的活性化スコアとの間のそれぞれの差を使用して前記訓練されていないモデルを訓練することであって、前記訓練すること(C)が、前記差に応答して、前記訓練されていないモデルと関連付けられた複数のパラメータを調整し、前記複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を前記目的の生理学的状態と関連付ける訓練されたモデルを得る、訓練することと、を含む、方法を実行するコンピュータによって実行可能な命令を集合的に符号化する、非一時的なコンピュータ可読媒体。

【請求項159】

前記モデルが、リグレッサーである、先行請求項のいずれか一項に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

本出願は、2021年6月15日に提出された「SYSTEMS AND METHODS FOR ASSOCIATING COMPOUNDS WITH PHYSIOLOGICAL CONDITIONS USING FINGERPRINT ANALYSIS」と題される米国仮特許出願第63/210,930号、及び2021年6月15日に提出された「COMPUTATIONAL MODELING PLATFORM」と題される63/210,679号の優先権を主張し、これらの各々は参照によりその全体が本明細書に組み込まれる。

【0002】

本発明は、概して、化合物を生理学的状態と関連付けるためのシステム及び方法に関する。

【背景技術】

【0003】

細胞メカニズムの研究は、疾患を理解するために重要である。

【0004】

生物組織は、動的かつ高度にネットワーク化された多細胞系である。特定の細胞における細胞内ネットワークの機能障害は、細胞行動の全体像を変化させ、疾患状態につながる。現在の創薬の努力は、細胞を健康な状態から疾患の状態へと遷移させる分子メカニズムを特徴付けることを目指し、これらの遷移を逆転又は阻害する薬理学的アプローチを特定する。これまでの努力はまた、これらの遷移を特徴付ける分子的特徴を特定し、これらの特徴を逆転させる薬理学的アプローチを特定することを目指していた。

【0005】

10

20

30

40

50

表面マーカーによって濃縮された組織又は細胞における細胞のバルク集合に関する分子データは、集団における個々の細胞の表現型及び分子多様性をマスクする。これらの細胞のバルク集合における細胞の不均一性は、疾患駆動メカニズムを解明することを目的とした現在の努力の結果を、誤解させるか、又は完全に不正確にさえさせる。単一細胞 RNA 配列決定などの新しいアプローチは、分子レベルで個々の細胞を特徴付けることができる。これらのデータは、より高い解像度で様々な細胞状態を理解するための基質を提供し、細胞が有する豊富で顕著な状態の多様性を明らかにする。

#### 【0006】

単一細胞データ、すなわち、これらのデータのまばらさ、細胞内に存在する分子の存在の見落とし、及びノイズを解釈する際に、これらの分子測定の精度に不確実性を伴う重大な課題が存在する。したがって、個々の細胞状態を制御するための薬理学的アプローチへの洞察を導き出し、それに応じて疾患を解決するために、新しいアプローチが必要である。

10

#### 【0007】

更に、複雑な疾患は、多くの場合、単一又はいくつかの分子標的に分解することができない。インビトロ疾患モデルのためのハイスループットイメージング技術及びハイスループットスクリーニングの最近の進歩にもかかわらず、インビトロベースのスクリーニングアプローチから生成された候補標的を有効な薬物に変換することは、多くの場合、比較的遅く、非効率的な分子標的ベースの創薬アプローチへの回帰を伴うかなりのタスクである。

20

#### 【0008】

上記の背景を考慮すると、当該技術分野で必要とされるのは、創薬のための候補化合物を識別するためのシステム及び方法である。

#### 【発明の概要】

#### 【0009】

本開示は、上記で特定された欠点に対処する。本開示は、少なくとも部分的に、目的の生理学的状態（例えば、表現型、疾患、細胞状態、及び/又は目的の細胞プロセス）に対応する細胞構成要素データ（例えば、遺伝子の存在量及び/又は振動シグネチャ）、並びに潜在表現及び機械学習を使用して、細胞構成要素のモジュール（例えば、サブセット）と、目的の生理学的状態との間の関連性（例えば、重み及び/又は相関）を決定することによって、これらの欠点に対処する。特に、本開示は、疾患などの様々な生理学的状態の基礎となる分子メカニズムを解明するためのシステム及び方法を提供する。

30

#### 【0010】

本開示の一態様は、試験化学化合物を目的の生理学的状態と関連付ける方法を提供する。この方法は、（A）試験化学化合物の化学構造のフィンガープリントを得ることを含む。

#### 【0011】

この方法は、（B）細胞構成要素モジュールのセットにアクセスすることを更に含む。細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々は、複数の細胞構成要素のそれぞれの独立したサブセットを含む。複数の細胞構成要素のそれぞれの独立したサブセットの各々についての対応する複数の細胞ベースのアッセイ存在量値は、生理学的状態と関連付けられた複数の異なる状態にわたって別々に相関する。細胞構成要素モジュールのセットにおける第1の細胞構成要素モジュールは、目的の生理学的状態と関連付けられる。

40

#### 【0012】

この方法は、（C）化学構造のフィンガープリントをモデルに入力することに対応して、モデルからの出力として、細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々についてのそれぞれの活性化スコアを取得することを更に含む。いくつかの実施形態において、モデルは、50以上のパラメータ、100以上のパラメータ、1000以上のパラメータ、又は10,000以上のパラメータを含む。

50

## 【 0 0 1 3 】

この方法は、(D)第1の細胞構成要素モジュールについての活性化スコアが、第1の閾値基準を満たす場合、試験化学化合物を目的の生理学的状態と関連付けることを更に含む。

## 【 0 0 1 4 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、臓器の細胞のものである。いくつかのそのような実施形態において、臓器は、心臓、肝臓、肺、筋肉、脳、膵臓、脾臓、腎臓、小腸、子宮、又は膀胱である。

## 【 0 0 1 5 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、組織の細胞のものである。いくつかの実施形態において、組織は、骨、軟骨、関節、気管、脊髄、角膜、眼、皮膚、又は血管である。 10

## 【 0 0 1 6 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、複数の幹細胞の細胞のものである。いくつかの実施形態において、複数の幹細胞は、複数の胚性幹細胞、複数の成体幹細胞、又は複数の人工多能性幹細胞(iPSC)である。

## 【 0 0 1 7 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、複数の初代ヒト細胞の細胞のものである。いくつかのそのような実施形態において、複数の初代ヒト細胞は、複数のCD34+細胞、複数のCD34+造血幹、複数の前駆細胞(HSPC)、複数のT細胞、複数の間葉系幹細胞(MSC)、複数の気道基底幹細胞、又は複数の人工多能性幹細胞である。 20

## 【 0 0 1 8 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、臍帯血中、末梢血中、又は骨髄中の細胞のものである。

## 【 0 0 1 9 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、固体組織中の細胞のものである。いくつかのそのような実施形態において、固体組織は、胎盤、肝臓、心臓、脳、腎臓、又は胃腸管である。

## 【 0 0 2 0 】

いくつかの実施形態において、細胞ベースのアッセイ存在量値は、複数の分化細胞のものである。いくつかのそのような実施形態において、複数の分化細胞は、複数の巨核球、複数の骨芽細胞、複数の軟骨細胞、複数の脂肪細胞、複数の肝細胞、複数の肝中皮細胞、複数の胆管上皮細胞、複数の肝星細胞、複数の肝類洞内皮細胞、複数のクッパー細胞、複数のピット細胞、複数の血管内皮細胞、複数の膵管上皮細胞、複数の膵管細胞、複数の腺房中心細胞、複数の腺房細胞、複数のランゲルハンス島、複数の心筋細胞、複数の線維芽細胞、複数のケラチノサイト、複数の平滑筋細胞、複数のI型肺胞上皮細胞、複数のII型肺胞上皮細胞、複数のクララ細胞、複数の線毛上皮細胞、複数の基底細胞、複数の杯細胞、複数の神経内分泌細胞、複数のクルチッキー(kultschitzky)細胞、複数の尿細管上皮細胞、複数の尿路上皮細胞、複数の円柱上皮細胞、複数の糸球体上皮細胞、複数の糸球体内皮細胞、複数の有足細胞、複数のメサングウム細胞、複数の神経細胞、複数の星状膠細胞、複数の小膠細胞、又は複数の乏突起膠細胞である。 30 40

## 【 0 0 2 1 】

いくつかの実施形態において、対応する複数の細胞ベースのアッセイ存在量値は、複数の細胞の単一細胞リボ核酸(RNA)配列決定(scRNA-seq)データである。いくつかのそのような実施形態において、生理学的状態に関連付けられた複数の異なる状態は、細胞のアリコートが生理学的状態に影響を与えることが知られている化合物に曝露されている対照状態に加えて、生理学的状態に影響を与えることが知られている1つ以上の参照化合物に異なる細胞のアリコートを曝露することによって導出される。

## 【 0 0 2 2 】



いくつかの実施形態において、対応する複数の細胞ベースのアッセイ存在量値は、バルクRNA配列に由来する。

【0023】

いくつかの実施形態において、対応する複数の細胞ベースのアッセイ存在量値は、単一細胞RNA配列決定に由来する。

【0024】

いくつかの実施形態において、細胞構成要素モジュールのセットは、第1の細胞構成要素モジュールからなる。

【0025】

いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールを含み、モデルは、複数のコンポーネントモデルを含むアンサンブルモデルである。複数のコンポーネントモデルにおけるコンポーネントモデルの各々は、化学構造のフィンガープリントを複数のコンポーネントモデルにおけるコンポーネントモデルの各々に入力することに対応して、細胞構成要素モジュールのセットにおける異なる細胞構成要素モジュールについての活性化スコアを提供する。

10

【0026】

いくつかの実施形態において、この方法は、試験化学化合物の単純化された分子入力ライブラリーシステム(SMILES)文字列表現からフィンガープリントを計算することを更に含む。

【0027】

いくつかのそのような実施形態において、複数のコンポーネントモデルにおけるコンポーネントモデルの各々は、対応するニューラルネットワーク(例えば、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせ)である。いくつかの実施形態において、対応するニューラルネットワークは、対応する完全に接続されたニューラルネットワーク及び対応するメッセージパッシングニューラルネットワークの組み合わせであり、対応する完全に接続されたニューラルネットワークの第1の出力及び対応するメッセージパッシングニューラルネットワークの第2の出力は、化学構造のフィンガープリントを対応する完全に接続されたニューラルネットワーク及び対応するメッセージパッシングニューラルネットワークに入力することに対応して組み合わせられ、細胞構成要素モジュールのセットにおける対応する細胞構成要素モジュールについての1つ以上の計算された活性化スコアにおける活性化スコアを決定する。

20

30

【0028】

いくつかのそのような実施形態において、複数のコンポーネントモデルにおけるコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルである。

【0029】

いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールであり、第1の細胞構成要素モジュールを含む複数の細胞構成要素モジュールの第1のサブセットは、目的の生理学的状態と関連付けられ、複数の細胞構成要素モジュールの第2のサブセットは、目的の生理学的状態と関連付けられず、第1の細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たし、複数の細胞構成要素モジュールの第2のサブセットにおける細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが、第1の閾値基準以外の第2の閾値基準を満たす場合、試験化学化合物は、目的の生理学的状態と識別される。

40

【0030】

いくつかの実施形態において、この方法は、電子形式で1つ以上の第1のデータセットを得、1つ以上の第1のデータセットが、第1の複数の細胞におけるそれぞれの細胞の各々について、第1の複数の細胞が、20個以上の細胞を含み、複数の注釈付きの細胞状態

50

を集合的に表し、複数の細胞構成要素（例えば、少なくとも10、20、30、100、若しくは1000個以上の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、それによって複数のベクトルにアクセスするか、又はそれらを形成することであって、複数のベクトルにおけるそれぞれのベクトルの各々が、(i) 複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii) 対応する複数のエレメントを含み、対応する複数のエレメントにおけるそれぞれのエレメントの各々が、第1の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を表す対応するカウントを有する、複数のベクトルにアクセスするか、又はそれらを形成すること、を含む、プロセスによって第1の細胞構成要素モジュールを識別することを更に含む。この方法は、複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することであって、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々が、複数の細胞構成要素のサブセットを含み、複数の細胞構成要素モジュールが、(i) 複数の候補細胞構成要素モジュール及び(ii) 複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、複数の細胞構成要素モジュールが、10を超える細胞構成要素モジュールを含む、識別すること、を更に含む。この方法は、電子形式で1つ以上の第2のデータセットを得、1つ以上の第2のデータセットが、第2の複数の細胞におけるそれぞれの細胞の各々について、第2の複数の細胞が、20個以上の細胞を含み、目的の生理学的状態を通知する複数の共変量を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、それによって、(i) 第2の複数の細胞、及び(ii) 複数の細胞構成要素又はその表現によって次元決定された細胞構成要素カウントデータ構造を得る。この方法は、複数の細胞構成要素又はその表現を共通次元として使用して細胞構成要素カウントデータ構造及び潜在表現を組み合わせることによって活性化データ構造を形成することであって、活性化データ構造が、複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、第2の複数の細胞における細胞の各々について、それぞれの活性化重みを含む、形成すること、複数の共変量におけるそれぞれの共変量の各々について、(i) 共変量のフィンガープリントの候補細胞構成要素モデルへの入力時に、候補細胞構成要素モデルによって表される細胞構成要素モジュールの各々に対する計算された活性化と、(ii) 候補細胞構成要素モデルによって表される細胞構成要素モジュールの各々に対する実際の活性化との間の差を使用して、候補細胞構成要素モデルを訓練することであって、訓練することが、差に応答して、候補細胞構成要素モデルと関連付けられた複数の共変量パラメータを調整する、訓練すること、を更に含む。いくつかのそのような実施形態において、複数の共変量パラメータは、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの共変量の各々について、それぞれの共変量が、第2の複数の細胞にわたって、それぞれの細胞構成要素モジュールと相関するかどうかを示す対応するパラメータを含み、方法が、候補細胞構成要素モデルを訓練する際に複数の共変量パラメータを使用して、複数の候補細胞構成要素モジュールにおける第1の細胞構成要素モジュールを識別することを更に含む。いくつかのそのような実施形態において、この方法は、複数の注射付きの細胞状態における注射付きの細胞状態が、曝露条件下での化合物への第1の複数の細胞における細胞の曝露（例えば、曝露期間、化合物の濃度、又は曝露期間及び化合物の濃度の組み合わせ）である、注射付きの細胞状態を更に含む。

#### 【0031】

いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定のmRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである。

#### 【0032】

いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特

10

20

30

40

50

定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせであり、第 1 又は第 2 の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量が、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動 ( F R E T ) 測定によって決定される。

【 0 0 3 3 】

いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせであり、第 1 又は第 2 の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量が、単一細胞リボ核酸 ( R N A ) 配列決定 ( s c R N A - s e q )、s c T a g - s e q、配列決定を  
10  
使用したトランスポザーゼ - アクセス可能なクロマチンのための単一細胞アッセイ ( s c A T A C - s e q )、C y T O F / S C o P、E - M S / A b s e q、m i R N A - s e q、C I T E - s e q、又はそれらの任意の組み合わせによって決定される。

【 0 0 3 4 】

いくつかの実施形態において、複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することは、複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、複数のベクトルに相関モデルを適用することを含む。いくつかのそのような実施形態において、相関モデルは、グラフクラスタリング (例えば、ピアソン相関ベースの距離メトリック上のライデン ( L e i d e n ) クラスタリング、ルーバン ( L o u v a i n ) クラスタリングなど)  
20  
を含む。

【 0 0 3 5 】

いくつかの実施形態において、複数の細胞構成要素モジュールは、10 ~ 2000 個の細胞構成要素モジュール、又は 100 ~ 8,000 個の細胞構成要素からなる。いくつかの実施形態において、複数の構成要素モジュールにおける候補細胞構成要素モジュールの各々は、200 ~ 300 個の細胞構成要素からなる。

【 0 0 3 6 】

いくつかの実施形態において、目的の生理学的状態は、疾患である。

【 0 0 3 7 】

いくつかの実施形態において、目的の生理学的状態は、疾患であり、第 1 の複数の細胞  
30  
が、複数の注釈付きの細胞状態によって示されるように、疾患を代表する細胞、及び疾患を代表しない細胞を含む。

【 0 0 3 8 】

いくつかの実施形態において、複数の共変量は、細胞バッチ、細胞ドナー、細胞型、疾患状態、化学化合物への曝露、又はそれらの任意の組み合わせを含む。

【 0 0 3 9 】

いくつかの実施形態において、候補細胞構成要素モデルを訓練することは、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する。  
40

【 0 0 4 0 】

いくつかの実施形態において、試験化学化合物は、2000 ダルトン未満の分子量を有する有機化合物である。いくつかのそのような実施形態において、試験化学化合物は、5 つの基準のリピンスキーの法則の各々を満たす有機化合物である。いくつかの実施形態において、試験化学化合物は、5 つの基準のリピンスキーの法則のうち少なくとも 3 つの基準を満たす有機化合物である。いくつかの実施形態において、モデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

【 0 0 4 1 】

10

20

30

40

50

いくつかの実施形態において、この方法は、Daylight、BCI、ECFP4、EcFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNN S2S、又はGraphConvを使用して、試験化学化合物の化学構造からフィンガープリントを生成することを更に含む。

【0042】

いくつかの実施形態において、細胞構成要素モジュールのセットは、5つ以上の細胞構成要素モジュール、10個以上の細胞構成要素モジュール、又は100個以上の細胞構成要素モジュールを含む。

【0043】

いくつかの実施形態において、それぞれの細胞構成要素モジュールにおける複数の細胞構成要素の独立したサブセットは、5つ以上の細胞構成要素を含む。 10

【0044】

いくつかの実施形態において、それぞれの細胞構成要素モジュールにおける複数の細胞構成要素の独立したサブセットは、目的の生理学的状態と関連付けられた分子経路における2~20個の細胞構成要素からなる。

【0045】

いくつかの実施形態において、第1の閾値基準は、第1の細胞構成要素モジュールが閾値活性化スコアを有することが必要である。

【0046】

本開示の別の態様は、試験化学化合物を目的の生理学的状態と関連付ける方法を提供する。 20

【0047】

この方法は、(A)試験化学化合物の化学構造のフィンガープリントを得ることを含む。

【0048】

この方法は、(B)摂動シグネチャのセットにアクセスすることであって、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、複数の細胞構成要素のそれぞれの独立したサブセットを含み、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である、アクセスすること、を更に含む。 30

【0049】

この方法は、(C)フィンガープリントをモデルに入力することであって、モデルが、50、100、500、1000、又は10,000以上のパラメータを含み、モデルが、フィンガープリントのモデルへの入力に回答して1つ以上の計算された活性化スコアを出力し、1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々が、摂動シグネチャのセットにおける対応する摂動シグネチャを表す、入力することを更に含む。 40

【0050】

この方法は、(D)摂動シグネチャのセットにおける第1の摂動シグネチャについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態と関連付けることを更に含む。

【0051】

いくつかの実施形態において、この方法は、試験化学化合物の単純化された分子入力ライントリーシステム(SMILES)文字列表現からフィンガープリントを計算することを更に含む。 50

## 【 0 0 5 2 】

いくつかの実施形態において、モデルは、ニューラルネットワークを含む。いくつかのそのような実施形態において、ニューラルネットワークは、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである。

## 【 0 0 5 3 】

いくつかの実施形態において、モデルは、複数のコンポーネントモデルを含むアンサンブルモデルであり、複数のコンポーネントモデルにおけるコンポーネントモデルの各々が、化学構造のフィンガープリントを複数のコンポーネントモデルのセットにおけるコンポーネントモデルの各々に入力することに応答して、摂動シグネチャのセットにおける異なる摂動シグネチャについての活性化スコアを提供する。

10

## 【 0 0 5 4 】

いくつかの実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

## 【 0 0 5 5 】

いくつかの実施形態において、複数のコンポーネントモデルにおけるコンポーネントモデルの各々は、対応するニューラルネットワークである（例えば、対応するニューラルネットワークは、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである）。

20

## 【 0 0 5 6 】

いくつかの実施形態において、複数のコンポーネントモデルにおけるコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルである。

## 【 0 0 5 7 】

いくつかの実施形態において、対応するニューラルネットワークは、完全に接続されたニューラルネットワーク及びメッセージパッシングニューラルネットワークの組み合わせであり、第1のニューラルネットワークの第1の出力及び第2のニューラルネットワークの第2の出力が、化学構造のフィンガープリントを完全に接続されたニューラルネットワーク及びメッセージパッシングニューラルネットワークに入力することに応答して、組み合わせられて、摂動シグネチャのセットにおける第1の摂動シグネチャについての1つ以上の計算された活性化スコアにおける活性化スコアを決定する。

30

## 【 0 0 5 8 】

いくつかの実施形態において、摂動シグネチャのセットは、複数の摂動シグネチャであり、第1の摂動シグネチャを含む、複数の摂動シグネチャの第1のサブセットが、目的の生理学的状態と関連付けられ、複数の摂動シグネチャの第2のサブセットが、目的の生理学的状態と関連付けられておらず、第1の摂動シグネチャについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たし、複数の摂動シグネチャの第2のサブセットにおける摂動シグネチャについてのそれぞれの計算された活性化スコアが、第1の閾値基準以外の第2の閾値基準を満たす場合、試験化学化合物が、目的の生理学的状態と識別される。

40

## 【 0 0 5 9 】

いくつかの実施形態において、目的の生理学的状態は、疾患である。

## 【 0 0 6 0 】

いくつかの実施形態において、試験化学化合物は、2000ダルトン未満の分子量を有する有機化合物である。

## 【 0 0 6 1 】

50

いくつかの実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則の各々を満たす有機化合物である。いくつかのそのような実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則のうち少なくとも3つの基準を満たす有機化合物である。

**【0062】**

いくつかの実施形態において、モデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

**【0063】**

いくつかの実施形態において、この方法は、Daylight、BCI、ECFP4、EcFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNN S2S、又はGraphConvを使用して、使用して、試験化学化合物の化学構造からフィンガープリントを生成することを更に含む。

**【0064】**

いくつかの実施形態において、摂動シグネチャのセットは、第1の摂動シグネチャからなる。

**【0065】**

いくつかの実施形態において、摂動シグネチャのセットは、5つ以上の摂動シグネチャ、10個以上の摂動シグネチャ、又は100個以上の摂動シグネチャを含む。

**【0066】**

いくつかの実施形態において、第1の閾値基準は、第1の摂動シグネチャが閾値活性化スコアを有することが必要である。

**【0067】**

本開示の別の態様は、化学化合物を目的の生理学的状態と関連付ける方法を提供する。

**【0068】**

この方法は、メモリ及び1つ以上のプロセッサを含むコンピュータシステムにおいて、(A)複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることを含む。

**【0069】**

この方法は、(B)複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む、得ること、を更に含む。

**【0070】**

この方法は、(C)複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、(ii)細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することであって、訓練すること(C)が、差に応答して訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、50、100、200、500、1000、又は10,000以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る、訓練すること、を更に含む。

**【0071】**

いくつかの実施形態において、細胞構成要素モジュールのセットは、単一の細胞構成要

10

20

30

40

50

素モジュールからなる。

【0072】

いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールを含む。

【0073】

いくつかの実施形態において、細胞構成要素モジュールのセットは、200～500個の細胞構成要素モジュールからなる。

【0074】

いくつかの実施形態において、複数の化合物は、 $10 \sim 1 \times 10^6$ 個の化合物からなる。

【0075】

いくつかの実施形態において、複数の化合物は、100～100,000個の化合物からなる。

【0076】

いくつかの実施形態において、複数の化合物は、1000～100,000個の化合物からなる。

【0077】

いくつかの実施形態において、訓練すること(C)は、回帰アルゴリズムに従って、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々についてのそれぞれの化合物の各々と関連付けられた差の各々に応答して、訓練されていないモデルと関連付けられた複数のパラメータを調整する。いくつかのそのような実施形態において、回帰アルゴリズムは、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々についてのそれぞれの化合物の各々と関連付けられた差の各々の最小二乗誤差を最適化する。

【0078】

いくつかの実施形態において、訓練されたモデルは、ニューラルネットワーク(例えば、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせ)を含む。

【0079】

いくつかの実施形態において、訓練されたモデルは、複数のコンポーネントモデルのアンサンブルモデルであり、複数のコンポーネントモデルにおけるそれぞれのコンポーネントモデルの各々が、複数の細胞構成要素モジュールにおける異なる細胞構成要素モジュールについて計算された活性化スコアを出力する。いくつかのそのような実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

【0080】

いくつかの実施形態において、複数のコンポーネントモデルにおけるコンポーネントモデルの各々は、対応するニューラルネットワークである。いくつかのそのような実施形態において、対応するニューラルネットワークは、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせである。

【0081】

いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールであり、複数の細胞構成要素モジュールの第1のサブセットは、目的の生理学的状態と関連付けられ、複数の細胞構成要素モジュールの第2のサブセットは、目的の生理学的状態と関連付けられていない。

【0082】

いくつかの実施形態において、この方法は、電子形式で1つ以上の第1のデータセットを得、1つ以上の第1のデータセットが、第1の複数の細胞におけるそれぞれの細胞の各

10

20

30

40

50

々について、第1の複数の細胞が、20個以上の細胞を含み、複数の注釈付きの細胞状態を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、複数の細胞構成要素が、5、10、15、20、25、50、又は100個以上の細胞構成要素を含み、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、それによって、複数のベクトルにアクセスするか、又はそれらを形成すること、を含む、プロセスによって複数の細胞構成要素モジュールにおける細胞構成要素モジュールを識別することを更に含む。複数のベクトルにおけるそれぞれのベクトルの各々は、(i)複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii)対応する複数のエレメントを含む。対応する複数のエレメントにおけるそれぞれのエレメントの各々は、第1の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を表す対応するカウントを有する。複数のベクトルは、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別するために使用され、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々は、複数の細胞構成要素のサブセットを含む。複数の細胞構成要素モジュールは、(i)複数の候補細胞構成要素モジュール、及び(ii)複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、複数の細胞構成要素モジュールは、3、5、10、15、20、又は100を超える細胞構成要素モジュールを含む。1つ以上の第2のデータセットは、電子形式で得られ、1つ以上の第2のデータセットは、第2の複数の細胞におけるそれぞれの細胞の各々について、第2の複数の細胞が、20個以上の細胞を含み、目的の生理学的状態を通知する複数の共変量を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含み、それによって、(i)第2の複数の細胞、及び(ii)複数の細胞構成要素又はその表現によって次元決定された細胞構成要素カウントデータ構造を得る。活性化データ構造は、複数の細胞構成要素又はその表現を共通次元として使用して、細胞構成要素カウントデータ構造及び潜在表現を組み合わせることによって形成され、活性化データ構造は、複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、第2の複数の細胞における細胞の各々について、それぞれの活性化重みを含む。候補細胞構成要素モデルは、(i)活性化データ構造を候補モデルに入力したときに、活性化データ構造内に表される細胞構成要素モジュールの各々における複数の共変量における各共変量の不在又は存在の予測と、(ii)細胞構成要素モジュールの各々における各共変量の実際の不在又は存在との間の差を使用して訓練される。この訓練は、差に応答して、候補細胞構成要素モデルと関連付けられた複数の共変量パラメータを調整する。

#### 【0083】

いくつかの実施形態において、複数の共変量パラメータは、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの共変量の各々について、それぞれの共変量が、第2の複数の細胞にわたって、それぞれの細胞構成要素モジュールと相関するかどうかを示す対応するパラメータを含み、候補細胞構成要素モデルを訓練する際に複数の共変量パラメータを使用して、複数の候補細胞構成要素モジュールにおける細胞構成要素モジュールを識別する。

#### 【0084】

いくつかの実施形態において、複数の注釈付きの細胞状態における注釈付きの細胞状態は、曝露条件下での化合物への第1の複数の細胞における細胞の曝露である。

#### 【0085】

いくつかの実施形態において、曝露条件は、曝露期間、化合物の濃度、又は曝露期間及び化合物の濃度の組み合わせである。

#### 【0086】

いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定のmRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである。



## 【0087】

いくつかの実施形態において、第1又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動(FRET)測定によって決定される。

## 【0088】

いくつかの実施形態において、第1又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、単一細胞リボ核酸(RNA)配列決定(scRNA-seq)、scTag-seq、配列決定を使用したトランスポーズ-アクセス可能なクロマチンのための単一細胞アッセイ(scATAC-seq)、CyTOF/SCoP、E-MS/Abseq、miRNA-seq、CITE-seq、又はそれらの任意の組み合わせによって決定される。

10

## 【0089】

いくつかの実施形態において、複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することは、複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、複数のベクトルに相関モデルを適用することを含む。いくつかのそのような実施形態において、相関モデルは、グラフクラスタリング(例えば、ピアソン相関ベースの距離メトリック上のライデン(Leiden)クラスタリング、又はルーバン(Louvain)クラスタリングである)を含む。

## 【0090】

いくつかの実施形態において、複数の細胞構成要素は、100~8,000個の細胞構成要素からなる。

20

## 【0091】

いくつかの実施形態において、複数の構成要素モジュールにおける候補細胞構成要素モジュールの各々は、200~300個の細胞構成要素からなる。

## 【0092】

いくつかの実施形態において、目的の生理学的状態は、疾患である。

## 【0093】

いくつかの実施形態において、生理学的状態は、疾患であり、第1の複数の細胞が、複数の注釈付きの細胞状態によって示されるように、疾患を代表する細胞、及び疾患を代表しない細胞を含む。

30

## 【0094】

いくつかの実施形態において、複数の共変量は、細胞バッチ、細胞ドナー、細胞型、疾患状態、又は化学化合物への曝露を含む。

## 【0095】

いくつかの実施形態において、候補細胞構成要素モデルを訓練することは、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する。

## 【0096】

いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、2000ダルトン未満の分子量を有する有機化合物である。

40

## 【0097】

いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則の各々を満たす。いくつかのそのような実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則のうち少なくとも3つの基準を満たす。

## 【0098】

いくつかの実施形態において、訓練されたモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最

50

近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

【0099】

いくつかの実施形態において、方法は、Daylight、BCI、ECFP4、ECFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNNS2S、又はGraphConvを使用して、対応する化学構造からそれぞれのフィンガープリントの各々を生成することを更に含む。

【0100】

いくつかの実施形態において、細胞構成要素モジュールのセットは、5つ以上の細胞構成要素モジュール、10個以上の細胞構成要素モジュール、又は100個以上の細胞構成要素モジュールを含む。

10

【0101】

本開示の別の態様は、化学化合物を目的の生理学的状態と関連付ける方法を提供する。この方法は、例えば、メモリ及び1つ以上のプロセッサを含むコンピュータシステムにおいて実施することができる。

【0102】

この方法は、(A)複数の化合物におけるそれぞれの化合物の各々の対応する化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることを含む。

【0103】

この方法は、(B)複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含む、得ることを更に含む。それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。

20

30

【0104】

この方法は、(C)複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、(ii)摂動シグネチャのセットにおける対応する化合物についてのそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。訓練(C)は、差に応答して、訓練されていないモデルと関連付けられた複数のパラメータを調整し、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。いくつかの実施形態において、複数のパラメータは、50、100、200、500、1000、10,000、又は $1 \times 10^6$ 以上のパラメータを含む。

40

【0105】

いくつかの実施形態において、摂動シグネチャのセットは、単一の摂動シグネチャからなる。

【0106】

いくつかの実施形態において、摂動シグネチャのセットは、200~500個の摂動シグネチャからなる。

【0107】

いくつかの実施形態において、複数の化合物は、10~ $1 \times 10^6$ 個の化合物からなる

50

。いくつかの実施形態において、複数の化合物は、100～100,000個の化合物からなる。いくつかの実施形態において、複数の化合物は、1000～100,000個の化合物からなる。

**【0108】**

いくつかの実施形態において、訓練すること(C)は、回帰アルゴリズムに従って、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々についての対応する化合物の各々と関連付けられた差の各々に応答して、訓練されていないモデルと関連付けられた複数のパラメータを調整する。いくつかのそのような実施形態において、回帰アルゴリズムは、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々についての対応する化合物の各々と関連付けられた差の各々の最小二乗誤差を最適化する。

10

**【0109】**

いくつかの実施形態において、訓練されたモデルは、ニューラルネットワーク(例えば、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせ)を含む。

**【0110】**

いくつかの実施形態において、訓練されたモデルは、複数のコンポーネントモデルのアンサンブルモデルであり、複数のコンポーネントモデルにおけるそれぞれのコンポーネントモデルの各々が、それぞれの化学構造のフィンガープリントを複数のコンポーネントモデルのセットにおけるコンポーネントモデルの各々に入力することに応答して、複数の摂動シグネチャのセットにおける異なる摂動シグネチャのセットについて計算された活性化スコアを出力する。いくつかのそのような実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

20

**【0111】**

いくつかの実施形態において、複数のコンポーネントモデルにおけるコンポーネントモデルの各々は、対応するニューラルネットワーク(例えば、完全に接続されたニューラルネットワーク、メッセージパッシングニューラルネットワーク、又はそれらの組み合わせ)である。

30

**【0112】**

いくつかの実施形態において、摂動シグネチャのセットは、複数の摂動シグネチャを含み、複数の摂動シグネチャの第1のサブセットは、目的の生理学的状態と関連付けられ、複数の摂動シグネチャの第2のサブセットは、目的の生理学的状態と関連付けられていない。

**【0113】**

いくつかの実施形態において、目的の生理学的状態は、疾患である。

**【0114】**

いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、2000ダルトン未満の分子量を有する有機化合物である。

40

**【0115】**

いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則の各々を満たす。

**【0116】**

いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則のうち少なくとも3つの基準を満たす。

**【0117】**

いくつかの実施形態において、訓練されたモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシン、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジステ

50

ニック回帰モデル、線形モデル、又は線形回帰モデルを含む。

【0118】

いくつかの実施形態において、方法は、Daylight、BCI、ECFP4、ECFC、MDL、APFP、TTFP、UNITY 2Dフィンガープリント、RNNS2S、又はGraphConvを使用して、対応する化学構造からそれぞれのフィンガープリントの各々を生成することを更に含む。

【0119】

いくつかの実施形態において、摂動シグネチャのセットは、5つ以上の摂動シグネチャ、10個以上の摂動シグネチャ、又は100個以上の摂動シグネチャを含む。

【0120】

いくつかの実施形態において、方法は、変化していない細胞状態と変化した細胞状態との間の差次的細胞構成要素存在量の尺度を表す単一細胞遷移シグネチャに電子形式でアクセスすることであって、変化した細胞状態が、変化していない細胞状態から変化した細胞状態への細胞遷移を通して発生し、(i)変化していない細胞状態、(ii)変化した細胞状態、及び(iii)変化していない細胞状態から変化した細胞状態への遷移のうちの少なくとも1つが、目的の生理学的状態と関連付けられ、単一細胞遷移シグネチャが、参照の複数の細胞構成要素の識別と、複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、変化していない細胞状態と変化した細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する第1の有意性スコアとを含む、アクセスすること、を含む、手順によって摂動シグネチャのセットにおけるそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアを得ることを更に含む。更に、単一細胞遷移シグネチャ及びそれぞれの摂動シグネチャを比較し、それによってそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定する。

【0121】

いくつかの実施形態において、単一細胞遷移シグネチャと摂動シグネチャとを比較して、それぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定することは、単一細胞遷移シグネチャの参照の複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの摂動シグネチャにおける対応する細胞構成要素の対応する有意性スコアに対するそれぞれの細胞構成要素の第1の有意性スコアを比較することを含む。

【0122】

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、摂動シグネチャのセットにおける他の摂動シグネチャと比較して、単一細胞遷移シグネチャに対するそれぞれの摂動シグネチャの関連性の相対的なランキングである。

【0123】

いくつかの実施形態において、相対的なランキングは、ウィルコクソンの順位和検定、t検定、ロジスティック回帰、又は一般化線形モデルによって決定される。

【0124】

いくつかの実施形態において、単一細胞遷移シグネチャの変化していない細胞状態が、それぞれの摂動シグネチャの第1の細胞状態又は第2の細胞状態と同じである。

【0125】

いくつかの実施形態において、単一細胞遷移シグネチャの変化していない細胞状態が、それぞれの摂動シグネチャの第1の細胞状態及び第2の細胞状態の両方とは異なる。

【0126】

いくつかの実施形態において、方法は、単一細胞遷移シグネチャの参照の複数の細胞構成要素、及びそれぞれの摂動シグネチャのそれぞれの複数の細胞構成要素を剪定して、転写因子と比較することを制限することを更に含む。

【0127】

いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動細胞状態は、複数の化合物における化合物に曝露されていない対照細胞によって表される。

10

20

30

40

50

## 【 0 1 2 8 】

いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動細胞状態は、それぞれの摂動シグネチャと関連付けられた化合物以外の複数の化学化合物における化学化合物に曝露されている無関係の摂動細胞にわたる平均によって表される。

## 【 0 1 2 9 】

開示された実施形態のいくつかにおいて、モデルは、リグレッサーである。

## 【 0 1 3 0 】

本開示の別の態様は、1つ以上のプロセッサ、及び1つ以上のプロセッサによる実行のための1つ以上のプログラムを格納するメモリを有するコンピュータシステムを提供し、1つ以上のプログラムは、本明細書に開示される方法及び/又は実施形態のうちのいずれかを実施するための命令を含む。

10

## 【 0 1 3 1 】

本開示の別の態様は、コンピュータによって実行するように構成された1つ以上のプログラムを格納する非一時的なコンピュータ可読記憶媒体を提供し、1つ以上のプログラムは、本明細書に開示される方法及び/又は実施形態のうちのいずれかを実行するための命令を含む。

## 【 0 1 3 2 】

本開示の更なる態様及び利点は、以下の詳細な説明から当業者に容易に明らかになるであろう。それにおいては、本開示の例示的な実施形態のみが示され、説明される。理解されるように、本開示は、他の異なる実施形態が可能であり、そのいくつかの詳細は、全て本開示から逸脱することなく、様々な明白な点で修正が可能である。したがって、図面及び説明は、本質的に例示とみなされるべきであり、限定とみなされるべきではない。

20

## 【 0 1 3 3 】

本明細書に開示される実施形態は、添付の図面の図において、限定としてではなく例として示される。同様の参照番号は、図面全体を通して対応する部分を指す。

## 【 図面の簡単な説明 】

## 【 0 1 3 4 】

【 図 1 】本開示の一実施形態による、例示的なシステム及びコンピューティングデバイスのブロック図を示す。

30

【 図 2 A 】本開示の様々な実施形態による、複数の細胞構成要素を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートをまとめて提供する。

【 図 2 B 】本開示の様々な実施形態による、複数の細胞構成要素を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートをまとめて提供する。

【 図 3 A 】本開示の様々な実施形態による、試験化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

【 図 3 B 】本開示の様々な実施形態による、試験化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

40

【 図 3 C 】本開示の様々な実施形態による、試験化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

【 図 3 D 】本開示の様々な実施形態による、試験化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

【 図 3 E 】本開示の様々な実施形態による、試験化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは

50

任意選択の要素を表す。

【図 4】本開示のいくつかの実施形態による、細胞構成要素の複数のベクトルの例及び細胞構成要素モジュールの潜在表現の例を示す。

【図 5】本開示のいくつかの実施形態による、細胞構成要素カウントデータ構造及び例示的な活性化データ構造の例を示す。

【図 6】本開示のいくつかの実施形態による、複数の化合物の重みを調整するためにモデルを訓練する方法の例を示す。

【図 7】本開示のいくつかの実施形態による、試験化学化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

10

【図 8】本開示の一実施形態による、化学化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

【図 9】本開示の一実施形態による、化学化合物を目的の生理学的状態と関連付けるための例示的な方法のプロセス及び特徴のフローチャートを提供し、破線のボックスは任意選択の要素を表す。

【図 10 A】本開示の一実施形態による、脂肪酸関連細胞プログラムの活性化のための化学構造を予測するための例示的な方法の性能及び 4 倍の検証を示す。図 10 A は、化学構造を予測するためのモデルアーキテクチャの概略図を示す。

【図 10 B】本開示の一実施形態による、脂肪酸関連細胞プログラムの活性化のための化学構造を予測するための例示的な方法の性能及び 4 倍の検証を示す。図 10 B は、1, 200 個のランダムに選択された化合物の試験セットにおける性能を示す。

20

【図 10 C】本開示の一実施形態による、脂肪酸関連細胞プログラムの活性化のための化学構造を予測するための例示的な方法の性能及び 4 倍の検証を示す。図 10 C は、訓練セットとは異なる足場を有する 1, 200 個の化合物の試験セットにおける性能を示す。

【図 10 D】本開示の一実施形態による、脂肪酸関連細胞プログラムの活性化のための化学構造を予測するための例示的な方法の性能及び 4 倍の検証を示す。図 10 D は、インビトロ前脂肪細胞アッセイにおける転写活性化に基づくベージング (beiging) 関連モジュールの検証を示す。

【図 10 E】本開示の一実施形態による、脂肪酸関連細胞プログラムの活性化のための化学構造を予測するための例示的な方法の性能及び 4 倍の検証を示す。図 10 E は、標的モジュールに対する 500 万個の化合物のデータベースから引き出された予測される化合物の最適化を示す。

30

【図 11】本開示の一実施形態による、胎児の赤血球生成及び T 細胞枯渇に関連する細胞挙動の活性化のための化学構造を予測するための例示的な方法の検証を示す。

【図 12】本開示の一実施形態による、単一細胞 RNA 配列決定 (scRNA-seq) を使用したヒト前脂肪細胞遺伝子モジュール活性化に対する既知のピペリジン含有化合物 (「KPC C」) 及び 6 つの新たに合成されたヒット「合成ヒット」の影響を評価するための例示的な方法の概略図を示す。

【図 13】本開示の一実施形態による、所望の転写変化の活性化に対する KPC C 及び 6 つの合成ヒットの効果を示す。

40

【図 14 A】任意選択の要素が破線のボックスによって示される細胞構成要素モジュールを識別するフローチャートを提供する。

【図 14 B】任意選択の要素が破線のボックスによって示される細胞構成要素モジュールを識別するフローチャートを提供する。

【図 14 C】任意選択の要素が破線のボックスによって示される細胞構成要素モジュールを識別するフローチャートを提供する。

【図 14 D】任意選択の要素が破線のボックスによって示される細胞構成要素モジュールを識別するフローチャートを提供する。

【発明を実施するための形態】

50

## 【0135】

導入。

上記の背景を考慮すると、本開示は、疾患に重要な細胞プロセス及びプログラムを標的とする創薬へのアプローチを記載する。このアプローチは、いくつかの態様において、生理学的状態（例えば、細胞プログラム、細胞プロセス、及び/又は細胞状態）及び化合物の化学構造のコンピュータにより操作された表現を使用して、化学構造関連モダリティ及びそれらの特性を予測することによって実現される。次いで、符号化された化学構造を細胞プログラム及び/又は細胞状態の表現にマッピングし、それによって、化合物を生理学的状態に関連付けることができる。

## 【0136】

例えば、いくつかの態様において、本開示は、分子プロファイル（例えば、遺伝子モジュール）と、目的の生物学的プロセス（例えば、細胞プログラム及び/又は細胞状態）及び化合物の化学構造との間の関連性を得るためのシステム及び方法を提供する。これらの関連性を使用して、創薬のために、類似の機能的又は構造的特性を有するものなどの新しい化学構造を予測することができる。

## 【0137】

いくつかの実施形態において、予測能力を有する計算モデリングアーキテクチャは、1つ以上のドメイン及び/又はデータタイプにわたる生理学的に関連する化学構造の潜在表現の生成を通じて、これらの関連性を発見するために使用される。関連性は、例えば、細胞の1つ以上の化合物への曝露に応答して、差次的遺伝子発現又は細胞状態遷移などの細胞挙動のプロファイルを提供する摂動データに由来し得る。いくつかの実施形態において、方法は、潜在表現及び機械学習を使用して、様々なドメイン（例えば、分子、細胞、臨床、インビボ、インビトロ、知識ベースなど）及び/又は様々なデータタイプ（転写、遺伝的、エピジェネティック、共変量など）の間の相関を組み合わせて決定して、生理学的に関連する化学構造を予測する。

## 【0138】

例示的な実施形態において、本開示は、化合物についての潜在表現を使用するモデリングアプローチを提供する。複数の化合物におけるそれぞれの化合物の各々について、方法は、それぞれの化合物が、複数の生理学的状態における生理学的状態の各々を誘発する可能性を表すベクトルを格納する潜在表現を生成することを含む。生理学的状態は、特定の表現型、細胞プロセス、及び/又は疾患と関連付けられた細胞状態遷移及び/又は細胞構成要素モジュール（例えば、遺伝子モジュール）を含むことができる。したがって、方法は、例えば、 $n$ \_\_化合物  $\times$   $n$ \_\_細胞\_\_状態又は  $n$ \_\_化合物  $\times$   $n$ \_\_遺伝子\_\_モジュールとして示される、化合物及び生理学的状態（例えば、細胞状態及び/又は遺伝子モジュール）によって次元決定されたモデルについてのマルチタスク訓練標識として機能するマトリックス表現を生成する。

## 【0139】

化合物を生理学的状態と関連付けるための機械学習モデルについての入力、化合物の化学構造を符号化し、更にモデルを訓練するために使用される、各化合物の正準異性体 SMILES 表現及び/又はグラフベースの表現を含む。訓練標識は、各化合物を各生理学的状態と関連付ける数値的活性化スコアとして提供される。例えば、各化合物についてのベクトルは、複数の関連する重みを含むことができ、各重みは、化合物が、それぞれの細胞状態、細胞状態遷移、摂動シグネチャ、及び/又はそれぞれの遺伝子モジュールの活性化などのそれぞれの生理学的状態を誘導する可能性を示す。

## 【0140】

入力としてマトリックス表現を受信すると、モデルは、回帰問題を解決することによって化学構造から細胞状態（例えば、摂動シグネチャ）及び/又は遺伝子モジュール活性化を学習するように訓練される。2つの例示的なモデルアーキテクチャは、回帰問題を解決するために使用される。第1のモデルは、SMILES文字列の標準的なフィンガープリント上で完全に接続されたネットワークを利用し、ネットワークアーキテクチャは、Re

10

20

30

40

50

LU活性化を伴う3層ネットワークである。第2のモデルは、DGLライブラリからのMPNNネットワークを含む。これらのモデルの各々は、回帰予測の最小二乗誤差を最適化することによって、互いに独立して訓練される。試験時間に、これらのモデルの予測は平均化され、したがって、第1及び第2のモデルを含むアンサンブルモデルを形成する。次いで、アンサンブルモデルを使用して、化合物と生理学的状態との間の関連性を決定することができ、これを更に適用して、化学構造から生理学的活性化の可能性の予測及び/又は特定の生理学的状態を誘発する可能性のある化学構造の予測を得ることができる。

#### 【0141】

有利には、本明細書に開示されるシステム及び方法は、創薬のための体系的でスケーラブルなアプローチを提供することによって、上記の欠点に対処する。例えば、創薬に関連する従来の機械学習アプローチは、ディープラーニング方法及び高性能コンピューティングとペアリングされた3Dタンパク質及び化学構造表現を使用したインシリコ標的スクリーニング能力を利用して、標的のライブラリに対する候補化合物の作用方法を計算する。しかしながら、これらのアプローチは、生物学的プロセスの基礎となる動的かつ高度にネットワーク化された多細胞系の複雑さに適切に対処していない、標的に焦点を当てたスクリーニングパラダイムに該当する。創薬のための他の従来の方法は、トランスクリプトームデータ又はイメージングデータに基づいて、単一の細胞及び細胞株が摂動にどのように応答するかをモデル化するために機械学習アプローチを使用する。そのような方法において、ハイスループットデータセットは、疾患の表現型表現及びインビトロ細胞系の複合摂動を学習するために使用される。これらは、表現型疾患応答を誘発又は相殺する化合物を予測するために使用される。しかしながら、従来のハイスループットデータモデリングアプローチは、それにもかかわらず、キュレーションの欠如及び多数の候補標的の識別の可能性によって不利になっている。ハイスループットスクリーニングから得られる潜在的な候補の各々の検証は、多くの場合、分子標的ベースの最適化又はインビトロスクリーニングのための数百若しくは更に数千又は化合物の合成を必要とする、手間のかかるプロセスである。

10

20

#### 【0142】

これらのアプローチとは対照的に、本開示は、次いで、生物学的プロセス（例えば、目的の生理学的状態に関与する遺伝子モジュール又は摂動シグネチャ）と関連付けられた細胞状態、摂動シグネチャ及び/又は細胞成分の表現にわたってマッピングされる、表現化学構造データ（例えば、化合物処理に対する細胞応答）を得るためのシステム及び方法を有利に提供する。それにもかかわらず、この標的に依存しないアプローチは、候補標的の体系的なキュレーション及び最適化を可能にし、したがって、標的発見とシステムにわたる予測翻訳との間のかなりのギャップを埋める。

30

#### 【0143】

例えば、以下の実施例に例示されるように、脂肪酸代謝に関与する候補ファーマコフォアは、本明細書に開示されるシステム及び方法の実施形態を使用して特定された。実施例4に更に例示されるように、候補ファーマコフォアに基づく予測翻訳は、6つの新しい化学物質を生成し、それらの全ては、ヒト脂肪細胞で試験した場合、脂肪酸関連細胞プロセスに関与する遺伝子モジュールを活性化することが見出された。タンパク質標的に対するハイスループットスクリーニング、特定若しくは最適化、又は数百若しくは数千の新しい化合物の合成を必要とせずに、候補ファーマコフォアの特定及び6つの新しい化学物質の設計を行った。したがって、本明細書に提供されるシステム及び方法は、標的発見から予測翻訳及び検証まで、従来の分子標的ベース又は表現型ベースのアプローチよりも、創薬及び開発プロセスの容易さ及び効率を改善する。

40

#### 【0144】

有利には、本開示は、化合物と生理学的状態との間の関連性（例えば、重み及び/又は相関）の標的化された決定のためのモデルの訓練及び使用を改善することによって、化合物と生理学的状態との関連性を改善する様々なシステム及び方法を更に提供する。機械学習モデルの複雑さは、時間の複雑性（所与の入力サイズnに対する実行時間、又はアルゴ

50



リズムの速度の尺度)、空間の複雑性(空間要件、又は所与の入力サイズ  $n$  に対するアルゴリズムを実行するために必要なコンピューティングパワー若しくはメモリの量)、又は両方を含む。複雑性(及びその後の計算負担)は、所与のモデルの訓練及び所与のモデルによる予測の両方に適用される。

#### 【0145】

いくつかの例では、計算の複雑性は、実装、追加のアルゴリズム若しくは交差検証方法の組み込み、及び/又は1つ以上のパラメータ(例えば、重み及び/又はハイパーパラメータ)によって影響を受ける。いくつかの例では、計算の複雑性は、入力サイズ  $n$  の関数として表され、入力データは、インスタンスの数(例えば、訓練試料の数)、次元  $p$ (例えば、特徴の数)、ツリー  $n_{trees}$  の数(例えば、ツリーに基づく方法の場合)、サポートベクトル  $n_{sv}$  の数(例えば、サポートベクトルに基づく方法の場合)、隣接  $k$  の数(例えば、 $k$  最近傍モデルの場合)、クラス  $c$  の数、及び/又は層  $i$  におけるニューロン  $n_i$  の数(例えば、ニューラルネットワークの場合)である。入力サイズ  $n$  に関して、次いで、(例えば、ビッグO表記での)計算の複雑性の近似は、入力サイズが増加するにつれて、実行時間及び/又は空間要件がどのように増加するかを示す。関数は、入力サイズの増加と比較して、より遅い速度又はより速い速度で複雑性を増加させることができる。計算の複雑性の様々な近似には、定数(例えば、 $O(1)$ )、対数(例えば、 $O(\log n)$ )、線形(例えば、 $O(n)$ )、対数線形(例えば、 $O(n \log n)$ )、二次(例えば、 $O(n^2)$ )、多項式(例えば、 $O(n^c)$ )、指数(例えば、 $O(c^n)$ )、及び/又は階乗(例えば、 $O(n!)$ )が含まれるが、これらに限定されない。いくつかの例では、定数関数の場合のように、入力サイズが増加するにつれて、より単純な関数はより低いレベルの計算の複雑性を伴うが、階乗関数などのより複雑な関数は、入力サイズのわずかな増加に回答して複雑性の大幅な増加を示すことができる。

#### 【0146】

機械学習モデルの計算の複雑性は、同様に(例えば、ビッグO表記で)関数によって表すことができ、複雑性は、モデルのタイプ、1つ以上の入力若しくは次元のサイズ、使用方法(例えば、訓練及び/若しくは予測)、並びに/又は時間若しくは空間の複雑性が評価されているかどうかに応じて変化し得る。例えば、決定木モデルにおける複雑性は、訓練のための  $O(n^2 p)$  及び予測のための  $O(p)$  として近似され、一方、線形回帰モデルにおける複雑性は、訓練のための  $O(p^2 n + p^3)$  及び予測のための  $O(p)$  として近似される。ランダムフォレストモデルの場合、訓練の複雑性は  $O(n^2 p n_{trees})$  として近似され、予測の複雑性は  $O(p n_{trees})$  として近似される。勾配ブーストモデルの場合、複雑性は、訓練のための  $O(n p n_{trees})$  及び予測のための  $O(p n_{trees})$  として近似される。カーネルサポートベクトルマシンの場合、複雑性は、訓練のための  $O(n^2 p + n^3)$ 、及び予測のための  $O(n_{sv} p)$  として近似される。ナイーブベイズモデルの場合、複雑性は、訓練のための  $O(np)$ 、及び予測のための  $O(p)$  として表され、ニューラルネットワークの場合、複雑性は、予測のための  $O(p n_1 + n_1 n_2 + \dots)$  として近似される。 $K$  最近傍モデルの複雑性は、時間のための  $O(k n p)$ 、及び空間のための  $O(n p)$  として近似される。ロジスティック回帰モデルの場合、複雑性は、時間のための  $O(n p)$ 、及び空間のため  $O(p)$  として近似される。ロジスティック回帰モデルの場合、複雑性は、時間のための  $O(n p)$ 、及び空間のため  $O(p)$  として近似される。

#### 【0147】

上述したように、機械学習モデルについて、計算の複雑性は、スケーラビリティを決定し、したがって、入力、特徴、及び/又はクラスサイズの増加、並びにモデルアーキテクチャのバリエーションのためのモデル(例えば、リグレッサー)の全体的な有効性及び有用性を決定する。大規模なデータセットの文脈において、少なくとも10、少なくとも100、少なくとも1000、又はそれ以上の細胞に対して得られた少なくとも10、少なくとも100、少なくとも1000、又はそれ以上の遺伝子の存在量を含む遺伝子発現データセットの場合と同様に、そのような大規模なデータセット上で実施される関数の計算

の複雑性は、多くの既存のシステムの能力に負担をかける可能性がある。更に、入力特徴の数（例えば、細胞構成要素（例えば、遺伝子）の数及び/又は化合物の数）及び/又はインスタンスの数（例えば、細胞の数、細胞状態注釈、摂動シグネチャ、モジュール、及び/又は共変量）が、技術的進歩とともに増加し、注釈の可用性を増加させ、下流の適用及び可能性を拡大するにつれて、任意の所与の分類モデルの計算の複雑性は、それぞれのシステムの仕様によって提供される時間及び空間容量を迅速に圧倒することができる。

#### 【0148】

したがって、化合物を生理学的状態と関連付けるための、最小入力サイズ（例えば、少なくとも10、少なくとも100、少なくとも1000、若しくはそれ以上の化合物；それぞれの細胞構成要素モジュールのための少なくとも10、少なくとも50、少なくとも100、若しくはそれ以上の細胞構成要素；少なくとも5、少なくとも10、少なくとも100、若しくはそれ以上の摂動シグネチャ；及び/又は少なくとも5、少なくとも10、少なくとも100、若しくはそれ以上の細胞構成要素モジュール）及び/又は対応する最小数のパラメータ（例えば、少なくとも50、少なくとも100、若しくは少なくとも1000のパラメータ及び/又は機械学習モデルに入力される特徴の全てのあらゆる可能なペアリングに対応するパラメータ）を有する機械学習モデルを使用することによって、計算の複雑性は、それが精神的に実施され得ないように比例して増加し、方法は、計算上の問題に対処する。例えば、本開示の一実施形態において、複数の少なくとも10個の細胞構成要素モジュール及び複数の少なくとも50個の化合物によって次元決定された活性化スコアマトリクスを得ることは、少なくとも500のパラメータ（例えば、重み）を得ることを含む。本開示の別の実施形態において、複数の少なくとも10個の摂動シグネチャにおける摂動シグネチャの各々について、複数の少なくとも50個の化合物における各化合物についてのそれぞれの活性化重みを得ることは、少なくとも500の活性化重みを得ることを含む。細胞状態遷移、細胞構成要素、細胞、化合物、共変量、試料、時点、複製、及び/又はバッチの数を含むがこれらに限定されない追加の入力特徴及び/又はインスタンスに同様の最小値を課すことは、同様に、方法の計算の複雑性に影響を与えるであろう。

#### 【0149】

機械学習モデルにおける計算の複雑性に関する更なる詳細は、2018年4月16日に公開され、[thekerneltrip.com/machine/learning/computational-complexity-learning-algorithms](http://thekerneltrip.com/machine/learning/computational-complexity-learning-algorithms)で、オンラインで入手可能な“Computational complexity of machine learning algorithms”、Hastie, 2001, The Elements of Statistical Learning, Springer, New York、及びArora and Barak, 2009, Computational Complexity: A Modern Approach, Cambridge University Press, New Yorkに提供されており、それらの各々は、参照によりその全体が本明細書に組み込まれる。

#### 【0150】

ここで、添付の図面にその例が示される実施形態を詳細に参照する。以下の詳細な説明には、本開示の完全な理解を提供するために、多くの特定の詳細が示されている。しかしながら、本開示は、これらの特定の詳細なしで実施され得ることが当業者に明らかである。他の場合では、実施形態の態様を不必要に曖昧にしないように、周知の方法、手順、構成要素、回路、及びネットワークは、詳細には説明されていない。

#### 【0151】

単一のインスタンスとして本明細書に記載される構成要素、動作、又は構造には、複数のインスタンスが提供され得る。最後に、様々な構成要素、動作、及びデータストアの間の境界は、多少任意であり、特定の動作は、特定の例示的な構成の文脈で例示される。他の形態の機能性が想定され、実装の範囲内に含まれ得る。概して、例示的な構成において

10

20

30

40

50

別個の構成要素として提示される構造及び機能性は、組み合わせられた構造又は構成要素として実装され得る。同様に、単一の構成要素として提示される構造及び機能性は、別個の構成要素として実装され得る。これら及び他の変形、修正、追加、及び改善は、実装の範囲内にある。

【0152】

また、「第1の」、「第2の」などの用語は、様々な要素を説明するために本明細書で使用され得るが、これらの要素はこれらの用語によって制限されるべきではないことも理解されたい。これらの用語は、ある要素を別の要素と区別するためにのみ使用される。例えば、第1のデータセットは、本発明の範囲から逸脱することなく、第2のデータセットと称され得、同様に、第2のデータセットは、第1のデータセットと称され得る。第1のデータセット及び第2のデータセットは両方ともデータセットであるが、同じデータセットではない。

10

【0153】

本明細書で使用される専門用語は、特定の実装を説明することのみを目的とし、特許請求の範囲を限定することを意図するものではない。実装及び添付の特許請求の範囲の説明で使用される場合、単数形の「a」、「an」及び「the」は、文脈により明らかにそうではないと指示されない限り、複数の形態も同様に含むことが意図される。「及び/又は (and/or)」という用語は、本明細書で使用するとき、関連する列挙した品目のうちの1つ以上の任意の及び全ての可能な組み合わせを指し、包含することも理解されるであろう。「含む (comprises)」及び/又は「含む (comprising)」という用語は、本明細書で使用される場合、記載された特徴、整数、ステップ、動作、要素、及び/又は成分の存在を指定するが、1つ以上の他の特徴、整数、ステップ、動作、要素、成分、及び/又はそれらの群の存在又は追加を排除しないことが更に理解されるであろう。

20

【0154】

本明細書で使用される場合、「～する場合 (if)」という用語は、文脈に応じて、「～するとき (when)」又は「～した後 (upon)」又は記載された先行する条件が真であることの「判定に応じて」又は「判定に従って」又は「検出に応じて」を意味すると解釈され得る。同様に、文脈に応じて、「(述べられた先行する条件が真であると) 判定される場合」又は「(述べられた先行する条件が真である) 場合」又は「(述べられた先行する条件が真である) とき」という語句は、述べられた先行する条件が真であることの「判定後」又は「判定に応じて」又は「判定に従って」又は「検出後」又は「検出に応じて」を意味すると解釈され得る。

30

【0155】

更に、参照番号が「i番目」の表示を与えられるとき、参照番号は、一般的な成分、セット、又は実施形態を指す。例えば、「細胞成分 i」と称される細胞成分は、複数の細胞成分における i 番目の細胞成分を指す。

【0156】

前述の説明は、例示的な実装を具現化する、例示的なシステム、方法、技術、命令シーケンス、及びコンピューティングマシンプログラム製品を含む。説明の目的において、本発明の主題の様々な実装の理解を提供するために、多くの特定の詳細が示されている。しかしながら、本発明の主題の実装は、これらの特定の詳細なしで実践され得ることは、当業者には明らかであろう。一般に、周知の命令インスタンス、プロトコル、構造、及び技術は、詳細に示されていない。

40

【0157】

説明の目的において、前述の説明は、特定の実装を参照して説明されている。しかしながら、以下の例示的な議論は、網羅的であることを意図するものではなく、又は、実装を開示される正確な形態に限定することを意図するものではない。上記の教示を考慮して、多くの修正及び変形が可能である。実装は、原理及びそれらの実際の用途を最もよく説明するために選択及び説明され、それによって、当業者が、企図される特定の使用に適した

50

実装及び様々な修正を伴う様々な実装を最もよく利用できるようにする。

【0158】

明確にするために、本明細書に記載される実装の慣例的特徴の全てが示され、説明されるわけではない。そのような任意の実際の実装の開発において、ユースケース及びビジネスに関連する制約への準拠など、設計者の特定の目標を達成するために多くの実装固有の決定が行われ、これらの特定の目標は、実装によって、及び設計者によって異なることが理解されるだろう。更に、そのような設計努力は複雑で時間がかかり得るが、それでも本開示の利益を得る当業者にとってはエンジニアリングの日常的な作業であることが理解されるであろう。

【0159】

本明細書のいくつかの部分は、情報に対する動作のアルゴリズム及び記号的表現の観点から、本発明の実施形態を説明する。これらのアルゴリズムの説明及び表現は、データ処理技術の当業者によって、それらの仕事の実質を当業者に効果的に伝達するために一般的に使用される。これらの動作は、機能的に、計算的に、又は論理的に説明されているが、コンピュータプログラム又は同等の電気回路、マイクロコードなどによって実装されることが理解される。

【0160】

本明細書で使用される言語は、可読性及び指示目的のために主に選択されており、本発明の主題を描写又は制限するために選択されていない場合がある。したがって、本発明の範囲は、この詳細な説明によって限定されるのではなく、それに基づく出願に関して生じる任意の特許請求の範囲によって限定されることが意図される。したがって、本発明の実施形態の開示は、本発明の範囲を例示することを意図するが、限定するものではない。

【0161】

一般に、特許請求の範囲及び本明細書で使用される用語は、当業者によって理解される平易な意味を有すると解釈されることが意図される。特定の用語は、追加の明確さを提供するために以下に定義される。明白な意味と提供される定義との間に矛盾がある場合、提供される定義が使用される。

【0162】

本明細書で直接定義されていない任意の用語は、本発明の技術分野内で理解されているように、それらに一般的に関連付けられた意味を有するものと理解されるべきである。ある特定の用語は、本発明の態様の組成物、デバイス、方法など、及びそれらを作製又は使用する方法を説明する際に、実践者に追加の指針を提供するために本明細書で議論される。同じことが複数の様式で言及され得ることが理解されるだろう。その結果、本明細書で議論される用語のうちの任意の1つ以上に対して、代替の言語及び同義語が使用され得る。用語が本明細書で詳述又は議論されるかどうかは重要ではない。いくつかの同義語又は置換可能な方法、材料などが提供される。1つ又はいくつかの同義語又は均等物の列挙は、それが明示的に述べられていない限り、他の同義語又は均等物の使用を排除しない。用語の例を含む例の使用は例示のみを目的とし、本明細書における本発明の態様の範囲及び意味を限定するものではない。

【0163】

定義。

本明細書で使用される場合、「約」又は「およそ」という用語は、当業者によって決定される特定の値に対する許容誤差範囲内であることを意味し、それは、部分的には、その値がどのように測定又は決定されるか、例えば、測定システムの限界に依存する。例えば、いくつかの実施形態において、「約」は、当該技術分野における慣例に従って、1以内又は1を超える標準偏差を意味する。いくつかの実施形態において、「約」は、所与の値の $\pm 20\%$ 、 $\pm 10\%$ 、 $\pm 5\%$ 、又は $\pm 1\%$ の範囲を意味する。いくつかの実施形態において、「約」又は「およそ」という用語は、値の1桁以内、5倍以内、又は2倍以内であることを意味する。本出願及び特許請求の範囲において特定の値が記載される場合、別段の記載がない限り、特定の値について許容可能な誤差範囲内であることを意味する「約」

10

20

30

40

50

という用語が、想定され得る。本明細書の詳細な説明内の全ての数値は、「約」示される値によって修正され、当業者によって予想される実験誤差及び変動を考慮する。「約」という用語は、当業者によって一般的に理解される意味を有することができる。いくつかの実施形態において、「約」という用語は、 $\pm 10\%$ を指す。いくつかの実施形態において、「約」という用語は、 $\pm 5\%$ を指す。

**【0164】**

本明細書で使用される場合、「存在量」、「存在量レベル」、又は「発現レベル」という用語は、1つ以上の細胞に存在する細胞構成要素（例えば、RNA種、例えば、mRNA若しくはmiRNA、又はタンパク質分子などの遺伝子産物）の量、又は複数の細胞にわたって存在する細胞構成要素の平均量を指す。mRNA又はタンパク質発現を指す場合、この用語は、一般に、特定の遺伝子座、例えば、特定の遺伝子に対応する任意のRNA又はタンパク質種の量を指す。しかしながら、いくつかの実施形態において、存在量は、複数のmRNA又はタンパク質アイソフォームを生じる特定の遺伝子に対応するmRNA又はタンパク質の特定のアイソフォームの量を指すことができる。遺伝子座は、遺伝子名、染色体位置、又は任意の他の遺伝子マッピングメトリックを使用して識別することができる。

10

**【0165】**

本明細書で同義的に使用される場合、「細胞状態」又は「生物学的状態」は、細胞又は細胞集団の状態又は表現型を指す。例えば、細胞状態は、健康であってもよい、又は疾患状態であってもよい。細胞状態は、複数の疾患のうちの一つであってもよい。細胞状態は、化合物治療及び/又は分化細胞系列に対する応答であってもよい。細胞状態は、1つ以上の遺伝子、1つ以上のタンパク質、及び/又は1つ以上の生物学的経路を含むが、これらに限定されない、1つ以上の細胞構成要素の尺度（例えば、活性化、発現、及び/又は存在量の尺度）によって特徴付けられ得る。

20

**【0166】**

本明細書で使用される場合、「細胞状態遷移」又は「細胞遷移」は、第1の細胞状態から第2の細胞状態への細胞の状態の遷移を指す。いくつかの実施形態において、第2の細胞状態は、変化した細胞状態（例えば、罹患した細胞状態への健康な細胞状態）である。いくつかの実施形態において、それぞれの第1の細胞状態及び第2の細胞状態のうち的一方は、非摂動状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方は、状態への細胞の曝露によって引き起こされる摂動状態である。摂動状態は、化合物への細胞の曝露によって引き起こされ得る。細胞状態遷移は、細胞内の細胞構成要素存在量の変化によって、したがって細胞（例えば、摂動シグネチャ）によって産生される同一性及び量の細胞構成要素（例えば、mRNA、転写因子）によってマークされ得る。

30

**【0167】**

本明細書で使用される場合、細胞又は複数の細胞についての細胞構成要素存在量測定に関連する「データセット」という用語は、いくつかの文脈において、単一細胞（例えば、単一細胞構成要素存在量データセット）から収集された高次元のデータセットを指すことができる。他の文脈では、「データセット」という用語は、単一細胞から収集された複数の高次元のデータセット（例えば、複数の単一細胞構成要素存在量データセット）を指すことができ、複数の細胞のうち一つの細胞から収集された複数のデータセットの各々を指すことができる。

40

**【0168】**

本明細書で使用される場合、「差次的存在量」又は「差次的発現」という用語は、第2の実体（例えば、第2の細胞、複数の細胞、及び/又は試料）と比較して、第1の実体（例えば、第1の細胞、複数の細胞、及び/又は試料）に存在する細胞構成要素の量及び/又は頻度の差を指す。いくつかの実施形態において、第1の実体は、第1の細胞状態（例えば、罹患した表現型）を特徴とする試料であり、第2の実体は、第2の細胞状態（例えば、正常又は健康な表現型）を特徴とする試料である。例えば、細胞構成要素は、第2の細胞状態を特徴とする実体と比較して、第1の細胞状態を特徴とする実体において高レベ

50

ル又は低レベルで存在するポリヌクレオチド（例えば、mRNA転写産物）であってもよい。いくつかの実施形態において、細胞構成要素は、第2の細胞状態を特徴とする実体と比較して、第1の細胞状態を特徴とする実体においてより高い頻度又はより低い頻度で検出されるポリヌクレオチドであってもよい。細胞構成要素は、量、頻度、又は両方の点で差次的に存在し得る。いくつかの場合において、一方の実体における細胞構成要素の量が、他方の実体における細胞構成要素の量と統計的に有意に異なる場合、細胞構成要素は、2つの実体の間で差次的に存在する。例えば、細胞構成要素は、他の実体に存在するものよりも、一方の実体において少なくとも約120%、少なくとも約130%、少なくとも約150%、少なくとも約180%、少なくとも約200%、少なくとも約300%、少なくとも約500%、少なくとも約700%、少なくとも約900%、若しくは少なくとも約1000%大きい場合、又は一方の実体において検出可能であり、他方の実体において検出不可能である場合、2つの実体において差次的に存在する。いくつかの場合において、実体の第1のサブセット（例えば、注釈付きの細胞状態の第1のサブセットを表す細胞）における細胞構成要素を検出する頻度が、実体の第2のサブセット（例えば、注釈付きの細胞状態の第2のサブセットを表す細胞）における頻度よりも統計的に有意に高い又は低い場合、細胞構成要素は、2つの実体のセットにおいて差次的に発現される。例えば、細胞構成要素は、一方の実体のセットにおいて、他の実体のセットよりも少なくとも約120%、少なくとも約130%、少なくとも約150%、少なくとも約180%、少なくとも約200%、少なくとも約300%、少なくとも約500%、少なくとも約700%、少なくとも約900%、又は少なくとも約1000%以上の頻度又は以下の頻度で観察される場合、2つの実体のセットにおいて差次的に発現される。

10

20

**【0169】**

本明細書で使用される場合、「健康な」という用語は、健康な状態（例えば、良好な健康を有する対象から得られた）を特徴とする試料を指す。健康な対象は、任意の悪性又は非悪性疾患の不在を示すことができる。「健康な」個体は、アッセイされる状態とは無関係であり、通常は「健康な」とみなすことができない他の疾患又は状態を有することができる。

**【0170】**

本明細書で使用される場合、細胞に関連する「摂動」という用語（例えば、細胞の摂動又は細胞摂動）は、1つ以上の化合物による治療などの1つ以上の状態への細胞の任意の曝露を指す。これらの化合物は、「ペルターバゲン（*perturbagens*）」と称され得る。いくつかの実施形態において、ペルターバゲンは、例えば、小分子、生物製剤、治療剤、タンパク質、小分子と組み合わせられたタンパク質、ADC、siRNA若しくは干渉RNAなどの核酸、cDNA過剰発現野生型及び/若しくは変異体shRNA、cDNA過剰発現野生型及び/若しくは変異体ガイドRNA（例えば、Cas9系若しくは他の遺伝子編集系）、又は前述のいずれかの任意の組み合わせを含むことができる。摂動は、細胞の表現型の変化、及び/又は細胞内の1つ以上の細胞構成要素の発現若しくは存在量レベルの変化（例えば、摂動シグネチャ）を誘発し得るか、又はそれによって特徴付けることができる。例えば、摂動は、細胞の転写プロファイルの変化によって特徴付けることができる。

30

40

**【0171】**

本明細書で使用される場合、「試料」、「生体試料」、又は「患者試料」という用語は、対象に関連する生物学的状態を反映し得る、対象から採取された任意の試料を指す。試料の例としては、対象の血液、全血、血漿、血清、尿、脳脊髄液、糞便、唾液、汗、涙、胸膜液、心膜液、又は腹膜液が挙げられるが、これらに限定されない。試料は、生きている又は死んでいる対象に由来する任意の組織又は材料を含むことができる。試料は、無細胞試料であってもよい。試料は、1つ以上の細胞構成要素を含むことができる。例えば、試料は、核酸（例えば、DNA若しくはRNA）若しくはその断片、又はタンパク質を含むことができる。「核酸」という用語は、デオキシリボ核酸（DNA）、リボ核酸（RNA）、又はそれらの任意のハイブリッド若しくは断片を指すことができる。試料中の核酸

50

は、無細胞核酸であってもよい。試料は、液体試料又は固体試料（例えば、細胞又は組織試料）であってもよい。試料は、体液であってもよい。試料は、糞便試料であってもよい。試料を処理して、組織又は細胞構造を物理的に破壊し（例えば、遠心分離及び/又は細胞溶解）、したがって、細胞内成分を、分析のために試料を調製するために使用され得る酵素、緩衝液、塩、洗剤などを更に含有し得る溶液中に放出することができる。

#### 【0172】

本明細書で使用される場合、化合物のフィンガープリントのような「フィンガープリント」という用語は、化合物のデジタルダイジェストである。そのようなデジタルダイジェストの非限定的な例としては、Daylightフィンガープリント、BCIFフィンガープリント、ECFC4フィンガープリント、ECFP4フィンガープリント、EcFCフィンガープリント、MDLフィンガープリント、原子対フィンガープリント（APFPフィンガープリント）、トポロジカル二面角フィンガープリント（TTFP）フィンガープリント、UNITY 2Dフィンガープリント、RNNS2Sフィンガープリント、又はGraphConvフィンガープリントが挙げられる。Franco, 2014, "The Use of 2D fingerprint methods to support the assessment of structural similarity in orphan drug legislation," J. Cheminform 6, p. 5、及びRensi and Altman, 2017, "Flexible Analog Search with Kernel PCA Embedded Molecule Vectors," Computational and Structural Biotechnology Journal, doi:10.1016/j.csbj.2017.03.003を参照されたく、それらの各々は参照により本明細書に組み込まれる。また、Raymond and Willett, 2002, "Effectiveness of graph-based and fingerprint-based similarity measures for virtual screening of 2D chemical structure databases," Journal of Computer-Aided Molecular Design 16, 59-71、及びFranco et al., 2014, "The use of 2D fingerprint methods to support the assessment of structural similarity in orphan drug legislation" Journal of chemoinformatics 6(5)も参照されたく、それらの各々は参照により本明細書に組み込まれる。

#### 【0173】

本明細書で使用される場合、「分類」という用語は、実体（例えば、細胞、試料、細胞構成要素、細胞構成要素モジュールなど）の特定の特性（例えば、細胞プロセス、共変量、細胞状態注釈など）に関連付けられた任意の数又は他の特徴を指すことができる。例えば、「+」記号（又は単語「正」）は、実体が特定の特性に対して正として分類されていることを示すことができる（例えば、細胞構成要素モジュールは、目的の細胞プロセスに正に関連付けられている）。別の例では、「分類」という用語は、実体と特定の特性との間の相関（例えば、それぞれの共変量とそれぞれの細胞構成要素モジュールとの間の相関）の決定を指すことができる。いくつかの実施形態において、分類は相関係数及び/又は重みである。分類は、バイナリ（例えば、正又は負）であってもよいが、又はより多くのレベルの分類（例えば、1~10又は0~1のスケール）を有してもよい。「カットオフ」及び「閾値」という用語は、動作で使用される所定の数を指すことができる。例えば、カットオフ値は、実体が除外される上記の値を参照することができる。閾値は、特定の分類が適用される値を上回るか、又は下回る値であってもよい。これらの用語のいずれかは、これらの文脈のいずれかで使用され得る。

#### 【0174】

本明細書で互換的に使用される場合、「分類子」、「モデル」、アルゴリズム、「リグ

10

20

30

40

50

レッサー」、及びノ」又は分類子」という用語は、機械学習モデル又はアルゴリズムを指す。いくつかの実施形態において、モデルは、教師なし学習アルゴリズムである。教師なし学習アルゴリズムの一例は、クラスタ分析である。

【0175】

いくつかの実施形態において、モデルは、教師あり機械学習である。教師あり学習アルゴリズムの非限定的な例としては、限定されないが、ロジスティック回帰、ニューラルネットワーク、サポートベクトルマシン、ナイーブベイズアルゴリズム、最近傍アルゴリズム、ランダムフォレストアルゴリズム、決定木アルゴリズム、ブーストツリーアルゴリズム、多項式ロジスティック回帰アルゴリズム、線形モデル、線形回帰、勾配ブースティング、混合モデル、隠れマルコフモデル、ガウシアンNBアルゴリズム、線形判別分析、又はそれらの任意の組み合わせが挙げられる。いくつかの実施形態において、モデルは、多項分類子アルゴリズムである。いくつかの実施形態において、モデルは、2段階の確率的勾配降下法(SGD)モデルである。いくつかの実施形態において、モデルは、ディープニューラルネットワーク(例えば、ディープアンドワイド試料レベルモデル)である。いくつかの実施形態において、本開示の分類子又はモデルは、25以上、100以上、1000以上、10,000以上、100,000以上、又は $1 \times 10^6$ 以上のパラメータを有するため、モデルの計算を精神的に実施することができない。

【0176】

更に、本明細書で使用される場合、「パラメータ」という用語は、アルゴリズム、モデル、リグレッサー、及びノ又は分類子における1つ以上の入力、出力、及びノ又は機能に影響を与える(例えば、修正、適応、及びノ又は調整する)ことができる、アルゴリズム、モデル、リグレッサー、及びノ又は分類子における内部又は外部エレメント(例えば、重み及びノ又はハイパーパラメータ)の任意の係数、又は同様に任意の値を指す。例えば、いくつかの実施形態において、パラメータは、アルゴリズム、モデル、リグレッサー、及びノ又は分類子の挙動、学習、及びノ又は性能を制御、修正、適応、及びノ又は調整するために使用され得る任意の係数、重み、及びノ又はハイパーパラメータを指す。いくつかの場合において、パラメータは、アルゴリズム、モデル、リグレッサー、及びノ又は分類子への入力(例えば、特徴)の影響を増加又は減少させるために使用される。非限定的な例として、いくつかの実施形態において、パラメータは、ノード(例えば、ニューラルネットワーク)の影響を増加又は減少させるために使用され、ノードは、1つ以上の活性化関数を含む。特定の入力、出力、及びノ又は関数へのパラメータの割り当ては、所与のアルゴリズム、モデル、リグレッサー、及びノ又は分類子のための任意の1つのパラダイムに限定されるものではなく、所望の性能のための任意の好適なアルゴリズム、モデル、リグレッサー、及びノ又は分類子アーキテクチャで使用することができる。いくつかの実施形態において、パラメータは、固定値を有する。いくつかの実施形態において、パラメータの値は、手動及びノ又は自動的に調整可能である。いくつかの実施形態において、パラメータの値は、アルゴリズム、モデル、リグレッサー、及びノ又は分類子のための検証及びノ又は訓練プロセスによって(例えば、誤差最小化及びノ又は逆伝搬方法によって)修正される。いくつかの実施形態において、本開示のアルゴリズム、モデル、リグレッサー、及びノ又は分類子は、複数のパラメータを含む。いくつかの実施形態において、複数のパラメータはn個のパラメータであり、ここで、 $n = 2; n = 5; n = 10; n = 25; n = 40; n = 50; n = 75; n = 100; n = 125; n = 150; n = 200; n = 225; n = 250; n = 350; n = 500; n = 600; n = 750; n = 1,000; n = 2,000; n = 4,000; n = 5,000; n = 7,500; n = 10,000; n = 20,000; n = 40,000; n = 75,000; n = 100,000; n = 200,000; n = 500,000; n = 1 \times 10^6; n = 5 \times 10^6$ 、又は $n = 1 \times 10^7$ である。したがって、本開示のアルゴリズム、モデル、リグレッサー、及びノ又は分類子は、精神的に実施することができない。いくつかの実施形態において、nは、 $10,000 \sim 1 \times 10^7$ 、 $100,000 \sim 5 \times 10^6$ 、又は $500,000 \sim 1 \times 10^6$ である。いくつかの実施形態において、本開示のアルゴリズム、モデル、リグレ

10

20

30

40

50



ッサー、及びノ又は分類子は、 $k$ 次元空間で動作し、ここで、 $k$ は、5又はそれよりも大きい（例えば、5、6、7、8、9、10など）正の整数である。したがって、本開示のアルゴリズム、モデル、リグレッサー、及びノ又は分類子は、精神的に実施することができない。

【0177】

ニューラルネットワーク。いくつかの実施形態において、モデルはニューラルネットワーク（例えば、畳み込みニューラルネットワーク及びノ又は残差ニューラルネットワーク）である。人工ニューラルネットワーク（ANN）としても知られるニューラルネットワークモデルは、畳み込み及びノ又は残差ニューラルネットワークモデル（ディープラーニングモデル）を含む。ニューラルネットワークは、入力データセットを出力データセットにマッピングするように訓練され得る機械学習モデルであり得、ニューラルネットワークは、ノードの複数の層に編成されたノードの相互接続されたグループを含む。例えば、ニューラルネットワークアーキテクチャは、少なくとも入力層、1つ以上の隠れ層、及び出力層を含み得る。ニューラルネットワークは、任意の総数の層、及び任意の数の隠れ層を含み得、隠れ層は、入力データのセットを出力値又は出力値のセットにマッピングすることを可能にする訓練可能な特徴抽出器として機能する。本明細書で使用される場合、ディープラーニングモデル（DNN）は、複数の隠れ層、例えば、2つ以上の隠れ層を含むニューラルネットワークであり得る。ニューラルネットワークの各層は、いくつかのノード（又は「ニューロン」）を含むことができる。ノードは、入力データ又は前の層のノードの出力のいずれかから直接来る入力を受信し、特定の動作、例えば、合計動作を実施することができる。いくつかの実施形態において、入力からノードへの接続は、パラメータ（例えば、重み及びノ又は重み係数）に関連付けられる。いくつかの実施形態において、ノードは、入力、 $x_i$ 、及びそれらに関連付けられたパラメータの全ての対の積を合計してもよい。いくつかの実施形態において、重み付けされた合計は、バイアス $b$ でオフセットされる。いくつかの実施形態において、ノード又はニューロンの出力は、線形関数又は非線形関数であってもよい閾値関数又は活性化関数 $f$ を使用してゲートされてもよい。活性化関数は、例えば、整流化線形ユニット（ReLU）活性化関数、漏洩ReLU活性化関数、又は飽和双曲線正接、同一性、バイナリストップ、ロジスティック、 $\arctan$ 、ソフトサイン、パラメトリック整流化線形ユニット、指数線形ユニット、 $\text{softplus}$ 、ベント同一性、 $\text{softexponential}$ 、正弦曲線、正弦、ガウシアン、若しくはシグモイド関数などの他の関数、又はそれらの任意の組み合わせであり得る。

【0178】

ニューラルネットワークの重み付け係数、バイアス値、及び閾値、又は他の計算パラメータは、訓練データの1つ以上のセットを使用して、訓練段階で「教示」又は「学習」され得る。例えば、パラメータは、ANNが計算する出力値が訓練データセットに含まれる例と一致するように、訓練データセットからの入力データ及び勾配降下又は後方伝搬法を使用して訓練され得る。パラメータは、逆伝搬ニューラルネットワーク訓練プロセスから取得され得る。

【0179】

様々なニューラルネットワークのいずれも、対象の画像を分析する際に使用するのに好適であり得る。例は、限定されないが、フィードフォワードニューラルネットワーク、放射基底関数ネットワーク、再帰ニューラルネットワーク、残差ニューラルネットワーク、畳み込みニューラルネットワーク、残差畳み込みニューラルネットワークなど、又はそれらの任意の組み合わせを含むことができる。いくつかの実施形態において、機械学習は、事前に訓練された及びノ若しくは転移学習されたANN又はディープラーニングアーキテクチャを利用する。畳み込み及びノ又は残差ニューラルネットワークは、本開示に従って対象の画像を分析するために使用することができる。

【0180】

例えば、ディープニューラルネットワークモデルは、入力層、複数の個別にパラメータ化された（例えば、重み付けされた）畳み込み層、及び出力スコアラーを含む。畳み込み

層の各々のパラメータ（例えば、重み）並びに入力層は、ディープニューラルネットワークモデルと関連付けられた複数のパラメータ（例えば、重み）に寄与する。いくつかの実施形態において、少なくとも100のパラメータ、少なくとも1000のパラメータ、少なくとも2000のパラメータ、又は少なくとも5000のパラメータは、ディープニューラルネットワークモデルに関連付けられる。そのため、ディープニューラルネットワークモデルは、精神的に解決され得ないため、コンピュータを使用する必要がある。換言すれば、モデルへの入力を与えられた場合、そのような実施形態において、モデル出力は、精神的にではなく、コンピュータを使用して決定される必要がある。例えば、Krizhevsky et al., 2012, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 2*, Pereira, Burges, Bottou, Weinberger, eds., pp. 1097 - 1105, Curran Associates, Inc., Zeiler, 2012 "ADADELTA: an adaptive learning rate method," 'CoRR, vol. abs/1212.5701、及びRumelhart et al., 1988, "Neurocomputing: Foundations of research," ch. Learning Representations by Back-propagating Errors, pp. 696 - 699, Cambridge, MA, USA: MIT Pressを参照されたく、それらの各々は参照により本明細書に組み込まれる。

10

20

#### 【0181】

モデルとしての使用に好適な畳み込みニューラルネットワークモデルを含む、ニューラルネットワークモデルは、例えば、Vincent et al., 2010, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J Mach Learn Res* 11, pp. 3371 - 3408、Larochelle et al., 2009, "Exploring strategies for training deep neural networks," *J Mach Learn Res* 10, pp. 1 - 40、及びHassoun, 1995, *Fundamentals of Artificial Neural Networks*, Massachusetts Institute of Technologyに開示されており、それらの各々は参照により本明細書に組み込まれる。モデルとしての使用に好適な更なる例示的なニューラルネットワークは、Duda et al., 2001, *Pattern Classification, Second Edition*, John Wiley & Sons, Inc., New York、及びHastie et al., 2001, *The Elements of Statistical Learning*, Springer-Verlag, New Yorkに開示されており、それらの各々は参照によりその全体が本明細書に組み込まれる。モデルとしての使用に好適な更なる例示的なニューラルネットワークはまた、Draghici, 2003, *Data Analysis Tools for DNA Microarrays*, Chapman & Hall/RC、及びMount, 2001, *Bioinformatics: sequence and genome analysis*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New Yorkにも記載されており、それらの各々は参照によりその全体が本明細書に組み込まれる。

30

40

#### 【0182】

サポートベクトルマシン。いくつかの実施形態において、モデルはサポートベクトルマシン(SVM)である。モデルとしての使用に好適なSVMモデルは、例えば、Cristianini and Shawe-Taylor, 2000, "An Introd

50

uction to Support Vector Machines," Cambridge University Press, Cambridge、Boser et al., 1992, "A training algorithm for optimal margin models," in Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory, ACM Press, Pittsburgh, Pa., pp. 142 - 152、Vapnik, 1998, Statistical Learning Theory, Wiley, New York、Mount, 2001, Bioinformatics: sequence and genome analysis, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y., Duda, Pattern Classification, Second Edition, 2001, John Wiley & Sons, Inc., pp. 259, 262 - 265、及び Hastie, 2001, The Elements of Statistical Learning, Springer, New York、及び Furey et al., 2000, Bioinformatics 16, 906 - 914 に記載されており、それらの各々は参照によりその全体が本明細書に組み込まれる。分類に使用される場合、SVMは、標識されたデータから、最大限に離れたハイパープレーンを使用して、バイナリ標識されたデータの所与のセットを分離する。線形分離が不可能な場合、SVMは、特徴空間への非線形マッピングを自動的に実現する「カーネル」の技術と組み合わせて機能することができる。特徴空間内のSVMによって見出されるハイパープレーンは、入力空間内の非線形決定境界に対応し得る。いくつかの実施形態において、SVMに関連付けられた複数のパラメータ（例えば、重み）は、ハイパープレーンを定義する。いくつかの実施形態において、ハイパープレーンは、少なくとも10、少なくとも20、少なくとも50、又は少なくとも100のパラメータによって定義され、SVMモデルは、それが精神的に解決され得ないため、計算するのにコンピュータを必要とする。

#### 【0183】

ナイーブベイズモデル。いくつかの実施形態において、モデルはナイーブベイズモデルである。モデルとしての使用に好適なナイーブベイズモデルは、例えば、Ng et al., 2002, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes," Advances in Neural Information Processing Systems, 14に開示されており、これは参照により本明細書に組み込まれる。ナイーブベイズ分類子は、特徴間の強い（ナイーブな）独立性仮定を伴うベイズの定理を適用することに基づく、「確率的分類子」のファミリー内の任意の分類子である。いくつかの実施形態において、それらは、カーネル密度推定と結合される。例えば、Hastie et al., 2001, The elements of statistical learning: data mining, inference, and prediction, eds. Tibshirani and Friedman, Springer, New Yorkを参照されたく、これは参照により本明細書に組み込まれる。

#### 【0184】

最近傍モデル。いくつかの実施形態において、モデルは、最近傍モデルである。最近傍モデルは、メモリベースであり得、適合するモデルを含まない。最近傍に関して、クエリ点  $x_0$ （試験対象）が与えられると、 $k$  個の訓練点  $x(r)$ 、 $r = 1, \dots, k$ （ここでは訓練対象）が、 $x_0$  に最も近い距離で識別され、次に点  $x_0$  が  $k$  個の最近傍を使用して分類される。ここで、これらの隣接までの距離は、識別遺伝子セットの存在量値の関数である。いくつかの実施形態において、特徴空間内のユークリッド距離は、 $d(i) = \|x(i) - x_0\|$  として距離を決定するために使用される。典型的には、最近傍モデルが使用される場合、線形判別を計算するために使用される存在量データは、平均ゼロ及

び分散 1 を有するように標準化される。最近傍法則は、不等クラス事前判定、差別的誤分類コスト、及び特徴選択の問題に対処するために改良することができる。これらの改良の多くは、隣接のための何らかの形態の重み付け投票を伴う。最近傍分析に関する更なる情報については、Duda, Pattern Classification, Second Edition, 2001, John Wiley & Sons, Inc、及び Hastie, 2001, The Elements of Statistical Learning, Springer, New York を参照されたく、それらの各々は参照により本明細書に組み込まれる。

#### 【0185】

k 最近傍モデルは、入力の特徴空間における k 個の最も近い訓練例からなる非パラメトリック機械学習方法である。出力はクラスメンバーシップである。オブジェクトは、隣接するオブジェクトの複数の投票によって分類され、オブジェクトは、その最も近い隣接する k 個の中で最も一般的なクラスに割り当てられる (k は、典型的には小さい正の整数である)。k = 1 の場合、オブジェクトは単にその単一の最も近い隣接のクラスに割り当てられる。参照により本明細書に組み込まれる、Duda et al., 2001, Pattern Classification, Second Edition, John Wiley & Sons を参照されたい。いくつかの実施形態において、k 最近傍モデルを解くために必要な距離計算の数は、それが精神的に実施され得ないために、コンピュータが所与の入力についてのモデルを解くために使用されるようなものである。

#### 【0186】

ランダムフォレスト、決定木、及びブーストツリーモデル。いくつかの実施形態において、モデルは、決定木である。モデルとしての使用に好適な決定木は、Duda, 2001, Pattern Classification, John Wiley & Sons, Inc., New York, pp. 395 - 396 に概説されており、これは参照により本明細書に組み込まれる。決定木に基づく方法は、特徴空間を長方形のセットに分割し、各々に (定数のような) モデルを適合させる。いくつかの実施形態において、決定木はランダムフォレスト回帰である。使用され得る 1 つの特定のモデルは、分類及び回帰木 (CART) である。他の特定の決定木モデルには、ID3、C4.5、MART、及びランダムフォレストが含まれるが、これらに限定されない。CART、ID3、及び C4.5 は、Duda, 2001, Pattern Classification, John Wiley & Sons, Inc., New York, pp. 396 - 408 and pp. 411 - 412 に記載されており、これは参照により本明細書に組み込まれる。CART、MART、及び C4.5 は、Hastie et al., 2001, The Elements of Statistical Learning, Springer-Verlag, New York, Chapter 9 に記載されており、これは参照によりその全体が本明細書に組み込まれる。ランダムフォレストは、Breiman, 1999, "Random Forests - Random Features," Technical Report 567, Statistics Department, U.C. Berkeley, September 1999 に記載されており、これは参照によりその全体が本明細書に組み込まれる。いくつかの実施形態において、決定木モデルは、少なくとも 10、少なくとも 20、少なくとも 50、又は少なくとも 100 のパラメータ (例えば、重み及び / 又は決定) を含み、それが精神的に解決され得ないため、計算するのにコンピュータを必要とする。

#### 【0187】

回帰。いくつかの実施形態において、モデルは回帰を使用する。回帰アルゴリズムは、任意のタイプの回帰であり得る。例えば、いくつかの実施形態において、回帰はロジスティック回帰である。いくつかの実施形態において、回帰は、ラッソ、L2 又は弾性ネット正規化によるロジスティック回帰である。いくつかの実施形態において、閾値を満たすことに失敗する対応する回帰係数を有するこれらの抽出された特徴は、考慮から取り除かれる (削除される)。いくつかの実施形態において、マルチカテゴリ応答を扱うロジスティ

ック回帰モデルの一般化が、モデルとして使用される。ロジスティック回帰は、Agresti, An Introduction to Categorical Data Analysis, 1996, Chapter 5, pp. 103 - 144, John Wiley & Son, New Yorkに開示されており、これは参照により本明細書に組み込まれる。いくつかの実施形態において、モデルは、Hastie et al., 2001, The Elements of Statistical Learning, Springer-Verlag, New Yorkに開示されている回帰モデルを利用する。いくつかの実施形態において、ロジスティック回帰モデルは、少なくとも10、少なくとも20、少なくとも50、少なくとも100、又は少なくとも1000のパラメータ（例えば、重み）を含み、それが精神的に解決され得ないため、計算するのにコンピュータを必要とする。 10

【0188】

線形判別分析。線形判別分析（LDA）、正規判別分析（NDA）、又は判別関数分析は、2つ以上のクラスの対象又はイベントを特徴付けるか、又は分離する特徴の線形組み合わせを見出すための統計学、パターン認識、及び機械学習で使用される方法であるフィッシャーの線形判別の一般化であり得る。得られる組み合わせは、本開示のいくつかの実施形態においてモデル（線形モデル）として使用され得る。

【0189】

混合モデル及び隠れマルコフモデル。いくつかの実施形態において、モデルは、McLachlan et al., Bioinformatics 18(3): 413 - 422, 2002に記載されるような混合モデルである。いくつかの実施形態において、特に、時間コンポーネントを含むそれらの実施形態において、モデルは、Schliep et al., 2003, Bioinformatics 19(1): i255 - i263に記載されるような隠れマルコフモデルである。 20

【0190】

クラスタリング。いくつかの実施形態において、モデルは教師なしクラスタリングモデルである。いくつかの実施形態において、モデルは教師ありクラスタリングモデルである。モデルとしての使用に好適なクラスタリングは、例えば、Duda and Hart, Pattern Classification and Scene Analysis, 1973, John Wiley & Sons, Inc., New York (本明細書以下では、“Duda 1973”)の211～256ページに記載されており、これは参照によりその全体が本明細書に組み込まれる。クラスタリング問題は、データセット内の自然なグルーピングを見出すことの1つとして記述することができる。自然なグルーピングを識別するために、2つの問題に対処することができる。第一に、2つの試料間の類似性（又は相違性）を測定する方法を決定することができる。このメトリック（例えば、類似性尺度）を使用して、1つのクラスタ内の試料が他のクラスタ内の試料よりも互いにより類似していることを確実にすることができる。第二に、類似性尺度を使用してデータをクラスタに分割するための機構を決定することができる。クラスタリング調査を開始する1つの方法は、距離関数を定義し、訓練セット内の試料の全てのペア間の距離のマトリックスを計算することであり得る。距離が類似性の良好な尺度である場合、同じクラスタ内の参照実体間の距離は、異なるクラスタ内の参照実体間の距離よりも有意に小さくてもよい。しかしながら、クラスタリングは、距離メトリックを使用しなくてもよい。例えば、ノンメトリック類似性関数  $s(x, x')$  を使用して、2つのベクトル  $x$  及び  $x'$  を比較することができる。 $s(x, x')$  は、 $x$  及び  $x'$  が何らかの形で「類似している」ときに値が大きい対称関数であり得る。データセット内の点間の「類似性」又は「相違性」を測定するための方法が選択されると、クラスタリングは、データの任意のパーティションのクラスタリング品質を測定する基準関数を使用することができる。基準関数を極端化するデータセットのパーティションを使用して、データをクラスタリングすることができる。本開示で使用することができる特定の例示的なクラスタリング技術は、階層的クラスタリング（最近傍アルゴリズム、最遠傍アルゴリズム、平均リンケージアルゴリズム、 30 40 50

重心アルゴリズム、又は二乗和アルゴリズムを使用した凝集クラスタリング)、k平均クラスタリング、ファジーk平均クラスタリング、及びジャービス・パトリック(Jarvis-Patrick)クラスタリングを含むことができるが、これらに限定されない。いくつかの実施形態において、クラスタリングは、教師なしクラスタリング(例えば、事前に考えられた数のクラスタ及び/又はクラスタ割り当ての事前決定を伴わない)を含む。

#### 【0191】

モデル及びブースティングのアンサンブル。いくつかの実施形態において、モデルのアンサンブル(2つ以上)が使用される。いくつかの実施形態において、AdaBoostなどのブースティング技術は、モデルの性能を改善するために、多くの他のタイプの学習アルゴリズムと併せて使用される。このアプローチでは、本明細書に開示されるモデルのいずれか、又はそれらの等価物の出力は、ブーストされたモデルの最終出力を表す加重合計に組み合わされる。いくつかの実施形態において、モデルからの複数の出力は、平均、中央値、モード、加重平均、加重中央値、加重モードなどを含むが、これらに限定されない、当該技術分野で既知の中心傾向の任意の尺度を使用して組み合わされる。いくつかの実施形態において、複数の出力は、投票方法を使用して組み合わされる。いくつかの実施形態において、モデルのアンサンブル内のそれぞれのモデルは、重み付けされるか、又は重み付けされない。

#### 【0192】

本明細書で使用される場合、「訓練されていないモデル」(例えば、「訓練されていないリグレッサー」及び/又は「訓練されていない分類子」という用語は、訓練データセットで訓練されていないリグレッサー又は分類子などの機械学習モデルを指す。本明細書で使用される場合、「モデルを訓練する」という用語は、訓練されていない、又は部分的に訓練されたモデルを訓練するプロセスを指す。例えば、いくつかの実施形態において、モデルを訓練することは、潜在表現で配置された複数の細胞構成要素モジュール及び以下で説明される細胞構成要素カウントデータ構造を得ることを含む。潜在表現及び細胞構成要素カウントデータ構造で配置された複数の細胞構成要素モジュールは、活性化データ構造(本明細書以下、「一次訓練データセット」)内の複数の細胞構成要素モジュールについての複数の共変量における各共変量の存在の実際の不在と併せて、訓練されていない又は部分的に訓練されたモデルに集合的な入力として適用される活性化データ構造を形成するために組み合わされて、共変量モジュール間で訓練されていない又は部分的に訓練されたモデルを訓練し、それによって訓練されたモデルを得る。更に、「訓練されていないモデル」という用語は、転移学習技術が訓練されていないモデルのそのような訓練に使用される可能性を排除しないことを理解されたい。例えば、参照により本明細書に組み込まれる、Fernandes et al., 2017, "Transfer Learning with Partial Observability Applied to Cervical Cancer Screening," Pattern Recognition and Image Analysis: 8th Iberian Conference Proceedings, 243 - 250は、そのような転移学習の非限定的な例を提供する。転移学習が使用される例では、上述の訓練されていないモデルは、一次訓練データセットのもの以上の追加データを提供される。すなわち、転移学習の実施形態の非限定的な例では、訓練されていないモデルは、(i)一次訓練データセット及び(ii)追加のデータを受信する。典型的には、この追加データは、別の補助訓練データセットから学習された係数(例えば、回帰係数)の形態である。更に、単一の補助訓練データセットの説明が開示されているが、本開示における訓練されていないモデルを訓練する際に一次訓練データセットを補完するために使用され得る補助訓練データセットの数に制限がないことを理解されたい。例えば、いくつかの実施形態において、2つ以上の補助訓練データセット、3つ以上の補助訓練データセット、4つ以上の補助訓練データセット、又は5つ以上の補助訓練データセットは、転移学習を通じて一次訓練データセットを補完するために使用され、そのような各補助データセットは一次訓練データセットと

10

20

30

40

50

は異なる。そのような実施形態において、転移学習の任意の方法が使用され得る。例えば、一次訓練データセットに加えて、第1の補助訓練データセット及び第2の補助訓練データセットが存在する場合を考慮する。第1の補助訓練データセットから（第1の補助訓練データセットへの回帰などのモデルの適用によって）学習された係数は、転移学習技術（例えば、2次元マトリクス乗算）を使用して第2の補助訓練データセットに適用され得、これは、次いで、係数が一次訓練データセットに適用される訓練された中間モデルをもたらし得、これは、一次訓練データセット自体と併せて、訓練されていないモデルに適用される。あるいは、第1の補助訓練データセットから学習された係数の第1のセット（第1の補助訓練データセットへの回帰などのモデルの適用によって）及び第2の補助訓練データセットから学習された係数の第2のセット（第2の補助訓練データセットへの回帰などのモデルの適用によって）は、各々個別に一次訓練データセットの別個のインスタンスに適用され得（例えば、別個の独立マトリクス乗算によって）、一次訓練データセット自体（又は一次訓練データセットから学習された主要コンポーネント若しくは回帰係数などの一次訓練データセットのいくつかの縮小形式）と併せて一次訓練データセットのインスタンスを分離するための係数のそのような適用の両方が、次いで、訓練されていないモデルを訓練するために訓練されていないモデルに適用され得る。いずれかの例では、第1及び第2の補助訓練データセットから導出される共変量モジュール相関（例えば、追加の細胞状態注釈、追加の共変量、及び/又はその細胞構成要素存在量など）に関する知識を、共変量で標識された一次訓練データセットと併せて、訓練されていないモデルを訓練するために使用される。

10

20

**【0193】**

本明細書で互換的に使用される場合、「ニューロン」、「ノード」、「ユニット」、「隠れニューロン」、「隠れユニット」などの用語は、入力を受け入れ、活性化関数及び1つ以上のパラメータ（例えば、係数及び/又は重み）を介して出力を提供するニューラルネットワークのユニットを指す。例えば、隠れニューロンは、以前の層からの1つ以上の入力を受け入れ、後続の層についての入力として機能する出力を提供することができる。いくつかの実施形態において、ニューラルネットワークは、1つの出力ニューロンのみを含む。いくつかの実施形態において、ニューラルネットワークは、複数の出力ニューロンを含む。一般的に、出力は、共変量、細胞状態注釈、又は目的の細胞プロセスなどの目的の状態の確率若しくは尤度、バイナリ判定（例えば、存在又は不在、正又は負の結果）、及び/又は標識（例えば、分類及び/又は相関係数）などの予測値である。単一クラス分類モデルの場合、出力は、状態（例えば、共変量、細胞状態注釈、及び/又は目的の細胞プロセス）を有する入力特徴（例えば、1つ以上の細胞構成要素モジュール）の尤度（例えば、相関係数及び/又は重み）であり得る。マルチクラス分類モデルの場合、複数の予測値を生成することができ、各予測値は、目的の状態の各々についての入力特徴の尤度を示す。

30

**【0194】**

本明細書で使用される場合、「パラメータ」という用語は、モデル、分類子、又はアルゴリズムにおける1つ以上の入力、出力、及び/又は機能に影響を与える（例えば、修正、適応、及び/又は調整する）ことができる、モデル、分類子、又はアルゴリズムにおける内部又は外部エレメント（例えば、重み及び/又はハイパーパラメータ）の任意の係数、又は同様に任意の値を指す。いくつかの実施形態において、パラメータは、モデルにおける1つ以上の入力、出力、又は関数を調節する係数（例えば、重み）である。例えば、パラメータの値を使用して、モデルへの入力（例えば、特徴）の影響をアップウェイト又はダウンウェイトすることができる。特徴は、ロジスティック回帰、SVM、又はナイーブベイズモデルなどのパラメータと関連付けることができる。パラメータの値は、代替的又は追加的に、ニューラルネットワークにおけるノード（例えば、ノードは、入力から出力への変換を定義する1つ以上の活性化関数を含む）、クラス、又はインスタンス（例えば、複数の細胞における細胞）の影響をアップウェイト又はダウンウェイトするために使用することができる。特定の入力、出力、機能、又は特徴へのパラメータの割り当ては、

40

50

所与のモデルのための任意の1つのパラダイムに限定されないが、最適な性能のための任意の適切なモデルアーキテクチャで使用され得る。いくつかの例では、モデルの入力、出力、機能、又は特徴と関連付けられたパラメータ（例えば、係数）への参照は、機械学習モデルの計算の複雑性のコンテキストなどにおいて、同じものの数、性能、又は最適化の指標として同様に使用され得る。いくつかの実施形態において、パラメータは、固定値を有する。いくつかの実施形態において、パラメータの値は、手動及び/又は自動的に（例えば、ハイパーパラメータ最適化方法を使用して）調整可能である。いくつかの実施形態において、パラメータの値は、モデル検証及び/又は訓練プロセスによって（例えば、本明細書の他の箇所に記載されるように、エラー最小化及び/又は逆伝搬方法によって）修正される。

10

**【0195】**

本明細書で使用される場合、「ベクトル」という用語は、エレメントの配列などのエレメントの列挙されたリストであり、各エレメントは割り当てられた意味を有する。したがって、本開示で使用される「ベクトル」という用語は、「テンソル」という用語と互換性がある。例として、ベクトルが存在量カウントを含む場合、複数の細胞において、それぞれの細胞構成要素について、複数の細胞の各々の1つについて、ベクトルに所定のエレメントが存在する。提示を容易にするために、いくつかの例では、ベクトルは、一次元であると説明され得る。しかしながら、本開示は、そのように限定されない。任意の次元のベクトルは、ベクトルにおける各エレメントが表すものの説明が定義されている（例えば、そのエレメント1は、複数の細胞の細胞1の存在量カウントなどを表す）ことを条件として、本開示で使用することができる。

20

**【0196】****I. 例示的なシステムの実施形態**

本開示のいくつかの態様の概要及び本開示で使用されるいくつかの定義が提供されたので、例示的なシステムの詳細は、図1と併せて説明される。

**【0197】**

図1は、本開示のいくつかの実施形態によるシステム100を示すブロック図を提供する。システム100は、目的の細胞プロセスと関連付けられた複数の細胞構成要素モジュールにおける1つ以上の細胞構成要素モジュールの決定を提供する。図1では、システム100はコンピューティングデバイスとして示されている。コンピュータシステム100の他のトポロジが可能である。例えば、いくつかの実施形態において、システム100は、実際には、ネットワーク内で一緒にリンクされるか、又はクラウドコンピューティング環境内で仮想マシン若しくはコンテナであるいくつかのコンピュータシステムを構成し得る。したがって、図1に示される例示的トポロジは、当業者に容易に理解されるような様式で、本開示の一実施形態の特徴を説明する役割を果たすだけである。

30

**【0198】**

図1を参照すると、いくつかの実施形態において、コンピュータシステム100（例えば、コンピューティングデバイス）は、ネットワークインターフェース104を含む。いくつかの実施形態において、ネットワークインターフェース104は、1つ以上の通信ネットワークを通じて（例えば、ネットワーク通信モジュール158を通じて）、システム内のシステム100コンピューティングデバイスを互いに、並びに任意選択の外部システム及びデバイスと相互接続する。いくつかの実施形態において、ネットワークインターフェース104は、インターネット、1つ以上のローカルエリアネットワーク（LAN）、1つ以上のワイドエリアネットワーク（WAN）、他のタイプのネットワーク、又はそのようなネットワークの組み合わせを介してネットワーク通信モジュール158を通じた通信を任意選択で提供する。

40

**【0199】**

ネットワークの例としては、ワールドワイドウェブ（WWW）、イントラネット及び/又は無線ネットワーク、例えば携帯電話ネットワーク、無線ローカルエリアネットワーク（LAN）及び/又は首都圏ネットワーク（MAN）、並びに無線通信による他のデバイ

50



スが挙げられる。無線通信は、グローバルモバイルコミュニケーションシステム（GSM）、エンハンスドデータGSM環境（EDGE）、高速ダウンリンクパケットアクセス（HSDPA）、高速アップリンクパケットアクセス（HSUPA）、エボリューション、データ専用（EV-DO）、HSPA、HSPA+、デュアルセルHSPA（DC-HSPDA）、ロングタームエボリューション（LTE）、近距離通信（NFC）、広帯域コード分割多重アクセス（W-CDMA）、コード分割多重アクセス（CDMA）、時分割多重アクセス（TDMA）、Bluetooth、ワイヤレスフィデリティ（Wi-Fi）（例えば、IEEE 802.11a、IEEE 802.11ac、IEEE 802.11ax、IEEE 802.11b、IEEE 802.11g及び/若しくはIEEE 802.11n）、ボイスオーバーインターネットプロトコル（VoIP）、Wi-MAX、電子メール用プロトコル（例えば、インターネットメッセージアクセスプロトコル（IMAP）及び/若しくはポストオフィスプロトコル（POP））、インスタントメッセージング（例えば、エクステンシブルメッセージング及びプレゼンスプロトコル（XMPP）、インスタントメッセージング及びプレゼンスレバレッジ拡張機能のセッション開始プロトコル（SIMPLE）、インスタントメッセージング及びプレゼンスサービス（IMPS））、並びに/又はショートメッセージングサービス（SMS）、あるいは本書の出願日の時点でまだ開発されていない通信プロトコルを含む任意の他の好適な通信プロトコルを含む、複数の通信規格、プロトコル及び技術のいずれかを任意選択で使用する。

10

#### 【0200】

20

いくつかの実施形態において、システム100は、1つ以上の処理ユニット（CPU）102（例えば、プロセッサ、処理コアなど）、1つ以上のネットワークインターフェース104、ユーザによって使用されるためのディスプレイ108及び入力システム105（例えば、入力/出力インターフェース、キーボード、マウスなど）を（任意選択で）含むユーザインターフェース106、メモリ（例えば、非永続的メモリ107、永続的メモリ109）、並びに前述のコンポーネントを相互接続するための1つ以上の通信バス103を含む。1つ以上の通信バス103は、システムコンポーネント間の通信を相互接続及び制御する回路（チップセットと呼ばれることもある）を任意選択で含む。非永続的メモリ107は、典型的には、DRAM、SRAM、DDR RAM、ROM、EEPROM、フラッシュメモリなどの高速ランダムアクセスメモリを含み、一方、永続的メモリ109は、典型的には、CD-ROM、デジタル汎用ディスク（DVD）、又は他の光学ストレージ、磁気カセット、磁気テープ、磁気ディスクストレージ、又は他の磁気記憶デバイス、磁気ディスク記憶デバイス、光学ディスク記憶デバイス、フラッシュメモリデバイス、又は他の不揮発性固体記憶デバイスを含む。永続的メモリ109は、任意選択で、CPU102から遠隔に位置する1つ以上の記憶デバイスを含む。永続的メモリ109、及び非永続的メモリ109内の不揮発性メモリデバイスは、非一時的コンピュータ可読記憶媒体を含む。いくつかの実施形態において、非永続的メモリ107又は代替的に、非一時的コンピュータ可読記憶媒体は、以下のプログラム、モジュール及びデータ構造、又はそれらのサブセットを、場合によっては永続的メモリ109と併せて格納する：

30

任意選択のオペレーティングシステム156（例えば、ANDROID、iOS、DARWIN、RTXC、LINUX、UNIX、OSX、WINDOWS、又はVxWorksなどの組み込みオペレーティングシステム）であって、様々な基本システムサービスを処理するための、及びハードウェア依存タスクを実施するための手順を含むオペレーティングシステム；

40

システム100を他のデバイス及び/又は通信ネットワーク104と接続するための任意選択のネットワーク通信モジュール（又は命令）158；

複数の化合物における化合物の各々についてのそれぞれの化学構造122（例えば、122-1、...、122-R）又はその表現（例えば、化学構造のフィンガープリント）を含む化合物構造データストア120；

細胞構成要素モジュール132のセット（例えば、132-1、...、132-K）

50

を含む細胞構成要素モジュールデータストア 130 であり、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素 134 のサブセット（例えば、134-1-1、...、134-1-Z）を含む；

摂動シグネチャ 142 のセット（例えば、142-1、...、142-P）を含む摂動データストア 140 であり、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別を含み、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第 1 の細胞状態とそれぞれの第 2 の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応するシグネチャスコア 144（例えば、144-1-1、...、144-1-Q）；

複数の化合物におけるそれぞれの化合物の各々について、それぞれの化学構造 152（例えば、152-1、...、152-R）の各々について、

任意選択で、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの数値的活性化スコア 154（例えば、154-1-1、...、154-1-K）、及び / 又は

任意選択で、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について、それぞれの数値的活性化スコア 156（例えば、156-1-1、...、156-1-P）を含む、活性化データ構造 150、並びに

複数のパラメータ（例えば、100 以上のパラメータ）を含むモデルであって、複数のパラメータは、計算された活性化スコアと、それぞれの化学構造についての数値的活性化スコアとの間の差に回答して調整される、モデル。

【0201】

様々な実施形態において、上述の識別されたエレメントのうちの 1 つ以上は、前述のメモリデバイスのうちの 1 つ以上に格納され、上述の機能を実施するための命令のセットに対応する。上記の識別されたモジュール、データ、又はプログラム（例えば、命令のセット）は、別個のソフトウェアプログラム、手順、データセット、又はモジュールとして実装される必要はなく、したがって、これらのモジュール及びデータの様々なサブセットは、様々な実装で組み合わせられてもよく、又は別様に再配置されてもよい。いくつかの実施形態において、非永続的メモリ 107 は、任意選択で、上記で識別されたモジュール及びデータ構造のサブセットを格納する。更に、いくつかの実施形態において、メモリは、上述されていない追加のモジュール及びデータ構造を格納する。いくつかの実施形態において、上記の識別されたエレメントのうちの 1 つ以上は、システム 100 のもの以外のコンピュータシステムに格納され、システム 100 によってアドレス指定可能であり、したがってシステム 100 は、必要なときにそのようなデータの全部又は一部を取り出すことができる。

【0202】

図 1 は、「システム 100」を描写するが、図は、本明細書に記載される実装の構造的な概略図ではなく、コンピュータシステムに存在し得る様々な特徴の機能的な説明としてより意図される。実際には、当業者によって認識されるように、別個に示される項目は組み合わせられてもよく、いくつかの項目は分離されてもよい。更に、図 1 は、非永続的メモリ 107 内のある特定のデータ及びモジュールを描写するが、これらのデータ及びモジュールの一部又は全ては、代わりに、永続的メモリ 109 内又は 2 つ以上のメモリ内に格納されてもよい。例えば、いくつかの実施形態において、少なくとも化合物構造データストア 120 及び活性化データ構造 150 は、クラウドベースのインフラストラクチャの一部であり得るリモート記憶デバイスに格納される。いくつかの実施形態において、少なくとも化合物構造データストア 120 及び活性化データ構造 150 は、クラウドベースのインフラストラクチャ上に格納される。いくつかの実施形態において、化合物構造データストア 120 及び活性化データ構造 150 もまた、リモート記憶デバイスに格納され得る。

【0203】

本開示によるシステムは、図 1 を参照して開示されているが、本開示による方法 200

10

20

30

40

50

、 300、700、800、900、及び1500は、図2、図3、図7、図8、図9、及び図14を参照してここで詳細に説明される。

【0204】

II. 試験化学化合物を目的の生理学的状態と関連付ける方法  
生理学的状態。

図3A～図3Eを参照すると、本開示の一態様は、試験化学化合物を目的の生理学的状態と関連付ける方法300を提供する。

【0205】

いくつかの実施形態において、目的の生理学的状態は、疾患である。

【0206】

いくつかの実施形態において、疾患は、感染性又は寄生虫性疾患、腫瘍、血液又は造血器官の疾患、免疫系の疾患、内分泌疾患、栄養疾患又は代謝疾患、精神障害、行動障害又は神経発達障害、睡眠覚醒障害、神経系の疾患、視覚系の疾患、耳又は乳様突起の疾患、循環器系の疾患、呼吸器系の疾患、消化器系の疾患、皮膚の疾患、筋骨格系又は結合組織の疾患、泌尿生殖器系の疾患、性的健康に関連する状態、妊娠、出産又は産褥期に関連する疾患、周産期に起因する特定の状態、及び発達異常からなる群から選択される。いくつかの実施形態において、疾患は、ICD-11 MMS、又は国際疾病分類の1つ以上の項目である。ICDは、疾患、負傷、及び死因を分類する方法を提供する。世界保健機関(WHO)は、診断された疾患の事例を記録及び追跡する方法を標準化するためにICDを発行している。

【0207】

いくつかの実施形態において、目的の生理学的状態は、疾患の前提条件又は併存疾患などの疾患刺激性である。

【0208】

いくつかの実施形態において、目的の生理学的状態は、細胞系で発生するか、又は細胞系の文脈で測定される。いくつかの実施形態において、目的の生理学的状態は、1つ以上の細胞において生じるか、又は1つ以上の細胞の文脈において測定され、1つ以上の細胞は、単一細胞、細胞株、生検試料細胞、及び/又は培養された初代細胞を含む。いくつかの実施形態において、目的の生理学的状態は、ヒト細胞において生じる生理学的状態である。いくつかの実施形態において、目的の生理学的状態は、本明細書に記載される試料(例えば、定義: 試料を参照されたい)のいずれかなどの試料において生じる生理学的状態である。いくつかの実施形態において、目的の生理学的状態は、ヒト又は動物などの対象において生じる生理学的状態である。

【0209】

いくつかの実施形態において、目的の生理学的状態は、目的の細胞プロセスであるか、又はそれに関連する。

【0210】

いくつかの実施形態において、目的の細胞プロセスは、異常な細胞プロセスである。いくつかの実施形態において、目的の細胞プロセスは、疾患と関連付けられた細胞プロセスである。例えば、上記のように、いくつかの実施形態において、この方法は、疾患に重要な細胞プロセス及びプログラムの標的化及び解明を提供する。いくつかの実施形態において、目的の細胞プロセスは、疾患の発症、進行、症状、重症度、及び/又は解消を含むが、これらに限定されない疾患の特徴のうちのいずれかの基礎となる機構を示すか、又はそれに関連する。いくつかの実施形態において、目的の細胞プロセスは、機能的経路である。いくつかの実施形態において、目的の細胞プロセスは、シグナル伝達経路である。いくつかの実施形態において、目的の細胞プロセスは、(例えば、化合物、小分子、及び/又は治療剤の)作用機序である。いくつかの実施形態において、目的の細胞プロセスは、転写ネットワーク(例えば、遺伝子調節ネットワーク)によって特徴付けられ、かつ/又は調節される。いくつかの実施形態において、目的の細胞プロセスは、第1の細胞状態と第2の細胞状態との間の遷移の間に生じる細胞プロセスである。

10

20

30

40

50

## 【0211】

いくつかの実施形態において、目的の細胞プロセスは、遺伝子セット濃縮アッセイ（GSEA）注釈、遺伝子オントロジー注釈、機能的及び/若しくはシグナル伝達経路注釈、並びに/又は細胞シグネチャ注釈などの注釈である。注釈は、NIH Gene Expression Omnibus（GEO）、EBI ArrayExpress、NCBI、BLAST、EMBL-EBI、GenBank、Ensembl、KEGG経路データベース、Library of Integrated Network-based Cellular Signatures（LINCS）L1000データセット、Reactome経路データベース、Gene Ontologyプロジェクト、及び/又は任意の疾患特異的データベースを含むが、これらに限定されない、任意の公知のデータベースから得ることができる。

10

## 【0212】

したがって、いくつかの実施形態において、目的の生理学的状態は、本明細書に記載される任意のそれぞれの疾患、機能的経路、シグナル伝達経路、作用機序、転写ネットワーク、不一致、及び/又は細胞若しくは生物学的プロセスである。

## 【0213】

いくつかの実施形態において、目的の生理学的状態は表現型である。例えば、いくつかの実施形態において、目的の生理学的状態は、化合物、小分子、及び/又は治療剤、例えば、疾患の毒性及び/又は解消などの生理学的兆候である。いくつかの実施形態において、生理学的状態は、フローサイトメトリーの読み出し、イメージング及び顕微鏡注釈（例えば、H&Eスライド、IHCスライド、放射線画像、及び/又は他の医学的イメージング）、並びに/又は細胞構成要素データを含むが、これらに限定されない実験データを使用して測定される表現型である。

20

## 【0214】

いくつかの実施形態において、目的の生理学的状態は、毒性の尺度である。いくつかの実施形態において、生理学的状態は、核受容体の阻害若しくは活性化、及び/又は核受容体の阻害の量若しくは活性化の量である。いくつかの実施形態において、生理学的状態は、阻害若しくは活性化、並びに/又は生物学的経路（例えば、ストレス応答経路）の阻害の量若しくは活性化の量である。本開示で使用され得る例示的な核受容体及び例示的なストレス応答経路、並びにこれらの核受容体及び例示的なストレス応答経路の阻害又は活性化データは、参照により本明細書に組み込まれるHuang et al., 2016, "Modelling the Tox21 10 K chemical profiles for in vivo toxicity prediction and mechanism characterization," Nat Commun. 7, p. 10425に記載のように、およそ10,000個の化合物に対して説明されている。

30

## 【0215】

いくつかの実施形態において、目的の生理学的状態は、細胞構成要素のセット（例えば、細胞構成要素モジュール）の活性化及び/又は摂動シグネチャ（例えば、摂動にตอบสนองする複数の分析物の差次的発現プロファイル）を特徴とする。

## 【0216】

例えば、いくつかの実施形態において、目的の生理学的状態は、細胞構成要素のセットを含む細胞構成要素モジュールである。任意の種類の実験物（例えば、遺伝子、転写物、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせ）は、それぞれの細胞構成要素モジュールにおける細胞構成要素のセットで使うことが企図される。いくつかの実施形態において、細胞構成要素モジュールは、当業者に明白であろうように、当業者に既知の任意の細胞又は生物学的プロセス、並びにその任意の異常と関連付けられる。本明細書に開示されるシステム及び方法とともに使用するのに好適な細胞構成要素モジュールは、以下の「細胞構成要素及び細胞構成要素モジュール」と題されるセクションに更に記載される。

40

## 【0217】

50

いくつかの実施形態において、目的の生理学的状態は、第1の細胞状態と第2の細胞状態との間の不一致を特徴とする摂動シグネチャ（例えば、細胞状態遷移シグネチャ）である。

**【0218】**

いくつかのそのような実施形態において、目的の生理学的状態は、疾患状態（例えば、疾患対象及び/又は疾患組織から得られた細胞）と健康な状態（例えば、健康又は対照の対象及び/又は組織から得られた細胞）との間の不一致によって識別される。例えば、いくつかの実施形態において、疾患状態は、細胞の機能の喪失、細胞の機能の獲得、細胞の進行（例えば、細胞の分化状態への遷移）、細胞の静止（例えば、細胞が分化状態に遷移することができない）、細胞の侵入（例えば、異常な位置における細胞の出現）、細胞の消失（例えば、細胞が通常存在する位置における細胞の不在）、細胞の障害（例えば、細胞内及び/又は周囲の構造、形態、及び/又は空間的变化）、細胞のネットワークの損失（例えば、子孫細胞又は細胞の下流の細胞における正常な効果を排除する細胞の変化）、細胞のネットワークの獲得（例えば、細胞の下流の細胞の子孫細胞における新しい下流の効果を引き起こす細胞の変化）、細胞の余剰（例えば、細胞の過剰）、細胞の不足（例えば、臨界閾値を下回る細胞の密度）、細胞内の細胞構成要素比及び/若しくは量の差、細胞における遷移速度の差、又はこれらの任意の組み合わせによって識別される。

10

**【0219】**

本明細書に開示されるシステム及び方法とともに使用するのに好適な摂動シグネチャは、以下の「摂動シグネチャ」と題されるセクションに更に記載される。

20

**【0220】**

いくつかの実施形態において、目的の生理学的状態は、複数の生理学的状態（例えば、細胞プロセス、細胞構成要素モジュール、及び/又は摂動シグネチャ）を含む。いくつかの実施形態において、目的の生理学的状態は、少なくとも3、少なくとも4、少なくとも5、少なくとも6、少なくとも7、少なくとも8、少なくとも9、少なくとも10、少なくとも15、少なくとも20、少なくとも30、少なくとも40、少なくとも50、少なくとも60、少なくとも70、少なくとも80、少なくとも90、又は少なくとも100の生理学的状態を含む。いくつかの実施形態において、目的の生理学的状態は、200以下、100以下、90以下、80以下、70以下、60以下、50以下、20以下、又は10以下の生理学的状態を含む。いくつかの実施形態において、目的の生理学的状態は、1~5、5~10、2~20、10~50、又は20~100の生理学的状態を含む。いくつかの実施形態において、目的の生理学的状態は、3以上の生理学的状態から始まり、200以下の生理学的状態で終わる別の範囲内にある複数の生理学的状態を含む。

30

**【0221】**

いくつかの実施形態において、本開示の化合物は、5つの基準のリピンスキーの法則を満たす化学化合物である。いくつかの実施形態において、本開示の化合物は、5つのリピンスキーの法則のうち2つ以上の法則、3つ以上の法則、又は4つ全ての法則を満たす有機化合物である。(i) 5つ以下の水素結合ドナー（例えば、OH及びNH基）、(ii) 10個以下の水素結合アクセプター（例えば、N及びO）、(iii) 500ダルトン未満の分子量、及び(iv) 5未満のLogP。4つの基準のうち3つが5という数字を含むため、「5つの法則」と呼ばれる。Lipinski, 1997, Adv. Drug Del. Rev. 23, 3を参照されたく、これは参照によりその全体が本明細書に組み込まれる。いくつかの実施形態において、本開示の化合物は、5つのリピンスキーの法則に加えて、1つ以上の基準を満たす。例えば、いくつかの実施形態において、本開示の化合物は、5個以下の芳香族環、4個以下の芳香族環、3個以下の芳香族環、又は2個以下の芳香族環を有する。

40

**【0222】**

ブロック302を参照すると、方法300は、試験化学化合物の化学構造のフィンガープリントを得ることを含む。

**【0223】**

50

例えば、いくつかの実施態様では、試験化学化合物を機械学習アプローチに適用することは、分子データ（例えば、化合物の化学構造）を機械学習モデルによって読み取り及び操作可能な形式に変換することを含む。

【0224】

図3Aのブロック304を参照すると、化学構造を機械学習読み取り可能なフォーマットに変換するための1つのアプローチは、テキストの文字列として分子を表す単純化された分子入力ラインエントリーシステム（SMILES）を使用して化学構造の「フィンガープリント」を決定することを含む。したがって、いくつかの実施形態において、この方法は、試験化学化合物の単純化された分子入力ラインエントリーシステム（SMILES）文字列表現からフィンガープリントを計算することを更に含む。SMILES文字列を使用した分子フィンガープリンティングは、例えば、Honda et al., 2019, "SMILES Transformer: Pre-trained Molecular Fingerprint for Low Data Drug Discovery," arXiv:1911.04738に更に記載されており、これは参照によりその全体が本明細書に組み込まれる。

10

【0225】

化学構造を機械学習読み取り可能フォーマットに変換するための別のアプローチは、グラフベースの分子フィンガープリントを決定することを含む。グラフベースの分子フィンガープリンティングでは、元の分子構造は、ノードが個々の原子を表し、エッジが原子間の結合を表すグラフによって表される。グラフベースのアプローチは、より低いサイズ要件を有する複数のサブ構造を効率的に符号化する能力、したがってより低い計算負荷、並びにフィンガープリント間の構造的類似性の表示を符号化する能力を含むいくつかの利点を提供する。グラフベースのフィンガープリンティングは、例えば、Duvenaud et al., 2015, "Convolutional networks on graphs for learning molecular fingerprints," NeurIPS, 2224-2232に更に記載され、これは参照によりその全体が本明細書に組み込まれる。いくつかの実施形態において、フィンガープリントはグラフ畳み込みネットワークから生成される。いくつかの実施形態において、フィンガープリントは、グラフアテンションネットワーク（GAT）、グラフ同形ネットワーク（GIN）、又はグラフ下部構造インデックススペース近似グラフ（SAGA）などの空間的グラフ畳み込みネットワークから生成される。いくつかの実施形態において、フィンガープリントは、チェビシェフ（Chebyshev）多項式フィルタリングを使用するスペクトルグラフ畳み込みなどのスペクトルグラフ畳み込みネットワークから生成される。

20

30

【0226】

図3Aのブロック306を参照すると、いくつかの実施形態において、フィンガープリントは、SMILES Transformer、ECFP4、RNNS2S、及び/又はGraphConvを使用して化学構造から生成される。

【0227】

モデルアーキテクチャ。

図3Bのブロック308を参照すると、方法は、フィンガープリントをモデルに入力することを含む。いくつかの実施形態において、モデルは、複数（例えば、100、200、300、500、1000、10,000又はそれ以上）のパラメータを含む。

40

【0228】

いくつかの実施形態において、モデルは、複数のパラメータ（例えば、重み及び/又はハイパーパラメータ）を含む。いくつかの実施形態において、モデルについての複数のパラメータは、少なくとも10、少なくとも50、少なくとも100、少なくとも500、少なくとも1000、少なくとも2000、少なくとも5000、少なくとも10,000、少なくとも20,000、少なくとも50,000、少なくとも100,000、少なくとも200,000、少なくとも500,000、少なくとも100万、少なくとも200万、少なくとも300万、少なくとも400万、又は少なくとも500万のパラメ

50

ータを含む。いくつかの実施形態において、モデルについての複数のパラメータは、800万以下、500万以下、400万以下、100万以下、500,000以下、100,000以下、50,000以下、10,000以下、5000以下、1000以下、又は500以下のパラメータを含む。いくつかの実施形態において、モデルについての複数のパラメータは、10~5000、500~10,000、10,000~500,000、20,000~100万、又は100万~500万のパラメータを含む。いくつかの実施形態において、モデルについての複数のパラメータは、10以上のパラメータから始まり、800万以下のパラメータで終わる別の範囲内にある。

#### 【0229】

いくつかの実施形態において、モデルの訓練は、1つ以上のハイパーパラメータ（例えば、訓練中に合わせられ得る1つ以上の値）によって更に特徴付けられる。いくつかの実施形態において、ハイパーパラメータ値は、訓練中に合わせられる（例えば、調整される）。いくつかの実施形態において、ハイパーパラメータ値は、訓練データセット及び/又は1つ以上の入力（例えば、細胞、細胞構成要素モジュール、共変量など）の特定のエレメントに基づいて決定される。いくつかの実施形態において、ハイパーパラメータ値は、実験的最適化を使用して決定される。いくつかの実施形態において、ハイパーパラメータ値は、ハイパーパラメータスイープを使用して決定される。いくつかの実施形態において、ハイパーパラメータ値は、以前のテンプレート又はデフォルト値に基づいて割り当てられる。

10

#### 【0230】

いくつかの実施形態において、1つ以上のハイパーパラメータのそれぞれのハイパーパラメータは、学習速度を含む。いくつかの実施形態において、学習速度は、少なくとも0.0001、少なくとも0.0005、少なくとも0.001、少なくとも0.005、少なくとも0.01、少なくとも0.05、少なくとも0.1、少なくとも0.2、少なくとも0.3、少なくとも0.4、少なくとも0.5、少なくとも0.6、少なくとも0.7、少なくとも0.8、少なくとも0.9、又は少なくとも1である。いくつかの実施形態において、学習速度は、1以下、0.9以下、0.8以下、0.7以下、0.6以下、0.5以下、0.4以下、0.3以下、0.2以下、0.1以下、0.05以下、0.01以下、又はそれ未満を含む。いくつかの実施形態において、学習速度は、0.0001~0.01、0.001~0.5、0.001~0.01、0.005~0.8、又は0.005~1である。いくつかの実施形態において、学習速度は、0.0001以上から始まり、1以下で終わる別の範囲内に入る。いくつかの実施形態において、1つ以上のハイパーパラメータは、正規化強度（例えば、L2重みペナルティ、中断率など）を更に含む。例えば、いくつかの実施形態において、モデル（例えば、ニューラルネットワーク）は、複数の隠れニューロンにおける隠れニューロンの各々の対応するパラメータ（例えば、重み）に対する正規化を使用して訓練される。いくつかの実施形態において、正規化は、L1又はL2ペナルティを含む。

20

30

#### 【0231】

いくつかの実施形態において、1つ以上のハイパーパラメータのそれぞれのハイパーパラメータは、損失関数である。いくつかの実施形態において、損失関数は、平均平方誤差、平滑化平均平方誤差、二次損失、平均絶対誤差、平均バイアス誤差、ヒンジ、マルチクラスサポートベクトルマシン、及び/又は交差エントロピーである。いくつかの実施形態において、損失関数は、勾配降下アルゴリズム及び/又は最小化関数である。

40

#### 【0232】

いくつかの実施形態において、モデルは、1つ以上の活性化関数と関連付けられる。いくつかの実施形態において、1つ以上の活性化関数における活性化関数は、tanh、シグモイド、softmax、ガウシアン、ボルツマン(Boltzmann)加重平均化、絶対値、線形、整流化線形ユニット(ReLU)、有界整流化線形、ソフト整流化線形、パラメータ化整流化線形、平均、最大、最初、サイン、平方、平方根、多重二乗、逆二乗、逆多重二乗、多調和スプライン、スイッチュ(swish)、ミッシュ(mish

50

)、ガウシアン誤差線形ユニット ( G e L U )、及び / 又は薄板スプラインである。モデルは、フィンガープリントのモデルへの入力に応答して1つ以上の計算された活性化スコアを出力する。

【 0 2 3 3 】

図 3 B のブロック 3 1 0 を参照すると、いくつかの実施形態において、モデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。いくつかの実施形態において、モデルは、リグレッサーである。いくつかの実施形態において、モデルは、本明細書に開示されるモデルのいずれかである ( 例えば、定義 : モデルを参照されたい ) 。

10

【 0 2 3 4 】

図 3 B のブロック 3 1 2 を参照すると、いくつかの実施形態において、モデルはニューラルネットワークを含む。

【 0 2 3 5 】

いくつかの実施形態において、ニューラルネットワークは、R e L U 活性化を有する完全に接続されたニューラルネットワークである。例えば、いくつかの実施形態において、モデルは、対応する1つ以上の入力を含むニューラルネットワークであり、対応する1つ以上の入力における各入力は、試験化学化合物についての化学構造のためのものであり、対応する複数の隠れニューロンを含む対応する第1の隠れ層であり、対応する複数の隠れニューロンにおける隠れニューロンの各々は、( i i ) 複数の入力における入力の各々に完全に接続され、( i i i ) 第1の活性化関数タイプと関連付けられ、( i i i i ) ニューラルネットワークについての複数のパラメータにおける対応するパラメータ ( 例えば、重み )、及び1つ以上の対応するニューラルネットワーク出力と関連付けられ、対応する1つ以上のニューラルネットワーク出力におけるそれぞれのニューラルネットワーク出力の各々は、( i i ) 入力として、対応する複数の隠れニューロンにおける隠れニューロンの各々の出力を直接的又は間接的に受信し、かつ( i i i ) 第2の活性化関数タイプと関連付けられる。いくつかのそのような実施形態において、ニューラルネットワークは、完全に接続されたネットワークである。

20

【 0 2 3 6 】

いくつかの実施形態において、ニューラルネットワークは、複数の隠れ層を含む。上述したように、隠れ層は、( 例えば、追加の複雑性を捕捉するために ) 入力層と出力層との間に位置する。複数の隠れ層が存在するいくつかの実施形態において、隠れ層の各々は、同じ又は異なるそれぞれの数のニューロンを有し得る。

30

【 0 2 3 7 】

いくつかの実施形態において、隠れニューロンの各々 ( 例えば、ニューラルネットワークにおけるそれぞれの隠れ層における ) は、入力データに対して関数 ( 例えば、線形関数又は非線形関数 ) を実施する活性化関数と関連付けられる。一般に、活性化関数の目的は、ニューラルネットワークが元のデータの表現について訓練され、その後、新しい ( 例えば、以前には見えなかった ) データの追加の表現を「適合」又は生成することができるように、データに非線形性を導入することである。特定の活性化関数は、データセットの極端な端部 ( 例えば、t a n h 及び / 又はシグモイド関数 ) で飽和をもたらす可能性があるため、活性化関数 ( 例えば、第1及び / 又は第2の活性化関数 ) の選択は、ニューラルネットワークの使用例に依存する。例えば、いくつかの実施形態において、活性化関数 ( 例えば、第1及び / 又は第2の活性化関数 ) は、本明細書に開示される任意の活性化関数を含むが、これらに限定されない、当該技術分野で既知の任意の好適な活性化関数から選択される。

40

【 0 2 3 8 】

いくつかの実施形態において、隠れニューロンの各々は、活性化関数に基づいて決定されたニューラルネットワークの出力に寄与するパラメータ ( 例えば、重み及び / 又はバイ

50



アス値)と更に関連付けられる。いくつかの実施形態において、隠れニューロンは、任意のパラメータ(例えば、ランダム化された重み)によって初期化される。いくつかの代替的な実施形態において、隠れニューロンは、所定のパラメータのセットによって初期化される。

#### 【0239】

いくつかの実施形態において、ニューラルネットワークにおける(例えば、1つ以上の隠れ層にわたる)複数の隠れニューロンは、少なくとも2個、少なくとも3個、少なくとも4個、少なくとも5個、少なくとも6個、少なくとも7個、少なくとも8個、少なくとも9個、少なくとも10個、少なくとも11個、少なくとも12個、少なくとも13個、少なくとも14個、少なくとも15個、少なくとも16個、少なくとも17個、少なくとも18個、少なくとも19個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、又は少なくとも500個のニューロンである。いくつかの実施形態において、複数の隠れニューロンは、少なくとも100個、少なくとも500個、少なくとも800個、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、少なくとも5000個、少なくとも6000個、少なくとも7000個、少なくとも8000個、少なくとも9000個、少なくとも10,000個、少なくとも15,000個、少なくとも20,000個、又は少なくとも30,000個のニューロンである。いくつかの実施形態において、複数の隠れニューロンは、30,000個以下、20,000個以下、15,000個以下、10,000個以下、9000個以下、8000個以下、7000個以下、6000個以下、5000個以下、4000個以下、3000個以下、2000個以下、1000個以下、又は500個以下のニューロンである。いくつかの実施形態において、複数の隠れニューロンは、2~20個、2~200個、2~1000個、10~50個、10~200個、20~500個、100~800個、50~1000個、500~2000個、1000~5000個、5000~10,000個、10,000~15,000個、15,000~20,000個、又は20,000~30,000個のニューロンである。いくつかの実施形態において、複数の隠れニューロンは、2個以上のニューロンから始まり、30,000個以下のニューロンで終わる別の範囲内にある。

10

20

30

#### 【0240】

いくつかの実施形態において、ニューラルネットワークは、1~50個の隠れ層を含む。いくつかの実施形態において、ニューラルネットワークは、1~20個の隠れ層を含む。いくつかの実施形態において、ニューラルネットワークは、少なくとも2個、少なくとも3個、少なくとも4個、少なくとも5個、少なくとも6個、少なくとも7個、少なくとも8個、少なくとも9個、少なくとも10個、少なくとも11個、少なくとも12個、少なくとも13個、少なくとも14個、少なくとも15個、少なくとも16個、少なくとも17個、少なくとも18個、少なくとも19個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、又は少なくとも100個の隠れ層を含む。いくつかの実施形態において、ニューラルネットワークは、100個以下、90個以下、80個以下、70個以下、60個以下、50個以下、40個以下、30個以下、20個以下、10個以下、9個以下、8個以下、7個以下、6個以下、又は5個以下の隠れ層を含む。いくつかの実施形態において、ニューラルネットワークは、1~5個、1~10個、1~20個、10~50個、2~80個、5~100個、10~100個、50~100個、又は3~30個の隠れ層を含む。いくつかの実施形態において、ニューラルネットワークは、1個以上の層から始まり、100個以下の層で終わる別の範囲内にある複数の隠れ層を含む。

40

#### 【0241】

いくつかの実施形態において、ニューラルネットワークは、浅いニューラルネットワー

50

クを含む。浅いニューラルネットワークは、少数の隠れ層を有するニューラルネットワークを指す。いくつかの実施形態において、そのようなニューラルネットワークアーキテクチャは、ニューラルネットワーク訓練の効率を改善し、訓練に關与する層の数の減少に起因して計算能力を節約する。いくつかの実施形態において、ニューラルネットワークは、1つの隠れ層を含む。いくつかの実施形態において、ニューラルネットワークは、2つ、3つ、4つ、又は5つの隠れ層を含む。

#### 【0242】

いくつかの実施形態において、ニューラルネットワークは、メッセージパッシングニューラルネットワークである。メッセージパッシングニューラルネットワークは、グラフ（例えば、化学構造のグラフベースの表現）上の教師あり学習のためのフレームワークを指し、ノードは原子を表し、エッジは原子間の結合を表す。一般に、メッセージパッシングニューラルネットワークは、フォワードパスにおける2つのフェーズ、メッセージパッシングフェーズ及び読み出しフェーズを含む。メッセージパッシングフェーズは、T間隔の期間にわたって実行され、メッセージ関数  $M_t$  及び頂点更新関数  $U_t$  に従って、グラフ内の各ノードで隠された状態を更新することを含む。読み出しフェーズは、読み出し関数  $R$  を使用してグラフについての特徴ベクトルを計算する。いくつかの実施形態において、メッセージパッシングニューラルネットワークは、畳み込みネットワーク（例えば、空間的グラフ畳み込みネットワーク及び/又はスペクトルグラフ畳み込みネットワーク）、ゲート付きグラフニューラルネットワーク（GG-NN）、相互作用ネットワーク、分子グラフ畳み込み、ディープテンソルニューラルネットワーク、及び/又はラブラシアンベースの方法を含む。例えば、Gilmer et al., 2017, "Neural Message Passing for Quantum Chemistry," arXiv:1704.01212v2を参照されたく、これは参照によりその全体が本明細書に組み込まれる。

10

20

#### 【0243】

図3Bのブロック314を参照すると、いくつかの実施形態において、モデルは、複数のコンポーネントモデルのアンサンブルモデルである。例えば、ブロック316を参照すると、いくつかの実施形態において、1つ以上の計算された活性化スコアにおける計算された活性化スコアの各々は、複数のコンポーネントモデルにおけるコンポーネントモデルの各々の出力の中心傾向の測定値である。

30

#### 【0244】

図3Bのブロック318を参照すると、いくつかの実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、線形回帰モデル、又は複数のニューラルネットワークを含む。

#### 【0245】

いくつかの実施形態において、アンサンブルモデルは、少なくとも2個、少なくとも3個、少なくとも4個、少なくとも5個、少なくとも6個、少なくとも7個、少なくとも8個、少なくとも9個、少なくとも10個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、又は少なくとも500個のコンポーネントモデルを含む。いくつかの実施形態において、アンサンブルモデルは、500個以下、400個以下、300個以下、200個以下、又は100個以下のコンポーネントモデルを含む。いくつかの実施形態において、アンサンブルモデルは、100個以下、50個以下、40個以下、30個以下、又は20個以下のコンポーネントモデルを含む。いくつかの実施形態において、アンサンブルモデルは、1~50個、2~20個、5~50個、10~80個、5~15個、3~30個、10~500個、2~100個、又は50~100個のコンポーネントモデルを含む。いくつかの実施形態において、アンサンブルモデルは、2個以上のコンポー

40

50

ネットモデルから始まり、500個以下のコンポーネントモデルで終わるコンポーネントモデルの別の範囲を含む。

【0246】

いくつかの実施形態において、アンサンブルモデルは、複数のコンポーネントモデルから得られた複数の出力（例えば、活性化スコア）を組み合わせることによって形成される。いくつかの実施形態において、分類子からの複数の出力（例えば、活性化スコア）は、平均、中央値、モード、加重平均、加重中央値、加重モード、算術平均、ミッドレンジ、ミッドヒンジ、トリミアン、及び/又はウィンザライズド平均を含むが、これらに限定されない、当該技術分野で既知の中心傾向の任意の尺度を使用して組み合わせられる。例えば、アンサンブルモデルからの最終決定は、アンサンブルモデル内の全てのコンポーネント

10

【0247】

いくつかの実施形態において、複数の出力は、投票方法を使用して組み合わせられる。例えば、いくつかの実施形態において、複数の出力は、それぞれの化学構造と、目的のそれぞれの生理学的状態との間の関連性を示す、アンサンブルモデル内のコンポーネントモデルの各々からの出力の数（例えば、活性化スコア）を集計することによって組み合わせられる。いくつかの実施形態において、コンポーネントモデルからの複数の出力（例えば、活性化スコア）は、多数決を使用して組み合わせられる。いくつかのそのような実施形態において、関連性を示す出力の集計（例えば、閾値基準を超える活性化スコアの集計）が投票閾値よりも大きい場合、それぞれの化学構造とそれぞれの目的の生理学的状態との間の関連性を決定することによって、コンポーネントモデルからの複数の出力が組み合わせられる。いくつかの実施形態において、投票閾値は、アンサンブルモデル内の複数のコンポーネントモデルからの総投票の少なくとも50%である。いくつかの実施形態において、投票閾値は、アンサンブルモデル内の複数のコンポーネントモデルからの総投票の少なくとも20%、少なくとも30%、少なくとも40%、少なくとも50%、少なくとも60%、少なくとも70%、少なくとも80%、少なくとも90%、又は少なくとも95%である。

20

【0248】

いくつかの実施形態において、アンサンブルモデル内のコンポーネントモデルの各々は、重み付けされない（例えば、各コンポーネントモデルは、アンサンブルモデル内で1票を有する）。いくつかの実施形態において、アンサンブルモデル内の1つ以上のコンポーネントモデルは、更に重み付けされる（例えば、アンサンブルモデル内で1票を超える投票を有する）。

30

【0249】

いくつかの実施形態において、方法は、単一のアンサンブルモデル又は複数のアンサンブルモデルを得ることを含む。当該技術分野で既知の任意のアーキテクチャが、アンサンブルモデルについて企図される。例えば、いくつかの実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、及び/又は線形回帰モデルを含む。いくつかの実施形態において、複数のコンポーネントモデルは、複数のニューラルネットワークを含む。

40

【0250】

図3Bのブロック320を参照すると、いくつかの実施形態において、モデルは、複数のニューラルネットワークのアンサンブルモデルである。図3Bのブロック322を参照すると、いくつかの実施形態において、モデルは、複数のニューラルネットワークを含むアンサンブルモデルであり、複数のニューラルネットワークにおける第1のニューラルネットワークは、ReLU活性化を伴う完全に接続されたニューラルネットワークであり、複数のニューラルネットワークにおける第2のニューラルネットワークは、メッセージパッシングニューラルネットワークである。いくつかのそのような実施形態において、第1

50

のニューラルネットワークは、入力として、化学構造についての分子フィンガープリントを SMILES 表現として受け入れる完全に接続された 3 層ニューラルネットワークである。いくつかの実施形態において、第 2 のニューラルネットワークは、入力として、化学構造についての分子フィンガープリントをグラフベースの表現として受け入れるメッセージパッシングニューラルネットワーク (MPNN) である。

#### 【0251】

細胞構成要素及び細胞構成要素モジュール。

上述のように、再びブロック 308 を参照すると、モデルへの化学構造についてのフィンガープリントの入力に応答して、モデルは、細胞構成要素モジュールのセットについて 1 つ以上の計算された活性化スコアを出力する。図 3C のブロック 326 を参照すると、1 つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、細胞構成要素モジュールのセットにおける対応する細胞構成要素モジュールを表す。

10

#### 【0252】

図 3C のブロック 328 を参照すると、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々は、複数の細胞構成要素の独立したサブセットを含む。

#### 【0253】

いくつかの実施形態において、細胞構成要素は、遺伝子、遺伝子産物 (例えば、mRNA 及び / 又はタンパク質)、炭水化物、脂質、エピジェネティック特徴、代謝産物、及び / 又はそれらの組み合わせである。いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである。いくつかの実施形態において、複数の細胞構成要素は、DNA、修飾 (例えば、メチル化) DNA、コード (例えば、mRNA) 又は非コード RNA (例えば、snRNA) を含む RNA を含む核酸、転写後修飾タンパク質 (例えば、リン酸化、グリコシル化、ミリスチル化などのタンパク質) を含むタンパク質、脂質、炭水化物、環状アデノシン一リン酸 (cAMP) 及び環状グアノシン一リン酸 (cGMP) などの環状ヌクレオチドを含む、ヌクレオチド (例えば、アデノシン三リン酸 (ATP)、アデノシン二リン酸 (ADP) 及びアデノシン一リン酸 (AMP))、酸化及び還元形態のニコチンアミドアデニンジヌクレオチド (NADP / NADPH) などの他の小分子細胞構成要素、並びにそれらの任意の組み合わせを含む。

20

30

#### 【0254】

いくつかの実施形態において、複数の細胞構成要素は、少なくとも 5 個、少なくとも 10 個、少なくとも 15 個、少なくとも 20 個、少なくとも 25 個、少なくとも 30 個、少なくとも 40 個、少なくとも 50 個、少なくとも 60 個、少なくとも 70 個、少なくとも 80 個、少なくとも 90 個、少なくとも 100 個、少なくとも 200 個、少なくとも 300 個、少なくとも 400 個、少なくとも 500 個、少なくとも 600 個、少なくとも 700 個、少なくとも 800 個、少なくとも 900 個、少なくとも 1000 個、少なくとも 2000 個、少なくとも 3000 個、少なくとも 4000 個、少なくとも 5000 個、少なくとも 6000 個、少なくとも 7000 個、少なくとも 8000 個、少なくとも 9000 個、少なくとも 10,000 個、少なくとも 20,000 個、少なくとも 30,000 個、少なくとも 50,000 個、又は 50,000 個を超える細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素は、70,000 個以下、50,000 個以下、30,000 個以下、10,000 個以下、5000 個以下、1000 個以下、500 個以下、200 個以下、100 個以下、90 個以下、80 個以下、70 個以下、60 個以下、50 個以下、又は 40 個以下の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素は、20 ~ 10,000 個の細胞構成要素からなる。いくつかの実施形態において、複数の細胞構成要素は、100 ~ 8,000 個の細胞構成要素からなる。いくつかの実施形態において、複数の細胞構成要素は、5 ~ 20 個、20 ~ 50 個、50 ~ 100 個、100 ~ 200 個、200 ~ 500 個、500 ~ 1000 個、1000 ~

40

50

5000個、5000～10,000個、又は10,000～50,000個の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素は、5個以上の細胞構成要素から始まり、70,000個以下の細胞構成要素で終わる別の範囲内にある。

【0255】

一例として、いくつかの実施形態において、複数の細胞構成要素は、RNAレベルで、任意選択で測定された複数の遺伝子を含む。いくつかの実施形態において、複数の遺伝子は、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも25個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも600個、少なくとも700個、少なくとも800個、少なくとも900個、又は少なくとも1000個の遺伝子を含む。いくつかの実施形態において、複数の遺伝子は、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、少なくとも5000個、少なくとも10,000個、少なくとも30,000個、少なくとも50,000個、又は50,000個を超える遺伝子を含む。いくつかの実施形態において、複数の遺伝子は、5～20個、20～50個、50～100個、100～200個、200～500個、500～1000個、1000～5000個、5000～10,000個、又は10,000～50,000個の遺伝子を含む。

10

【0256】

別の例として、いくつかの実施形態において、複数の細胞構成要素は、複数のタンパク質を含む。いくつかの実施形態において、複数のタンパク質は、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも25個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも600個、少なくとも700個、少なくとも800個、少なくとも900個、又は少なくとも1000個のタンパク質を含む。いくつかの実施形態において、複数のタンパク質は、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、少なくとも5000個、少なくとも10,000個、少なくとも30,000個、少なくとも50,000個、又は50,000個を超えるタンパク質を含む。いくつかの実施形態において、複数のタンパク質は、5～20個、20～50個、50～100個、100～200個、200～500個、500～1000個、1000～5000個、5000～10,000個、又は10,000～50,000個のタンパク質を含む。

20

30

【0257】

細胞構成要素モジュールにおける細胞構成要素の各々が一意であるという要件はない。例えば、細胞構成要素モジュールAが、細胞構成要素1、3及び10を含有する場合を考慮する。細胞構成要素モジュールのセットにおける他の細胞構成要素モジュールは、これらの細胞構成要素も含有し得る。ここで、「独立した」という用語は、特定の細胞構成要素モジュールにおける複数の細胞構成要素のサブセットが全体として一意であることを意味する。したがって、上記の例示的な細胞構成要素モジュールAを考慮すると、細胞構成要素モジュールのセットにおける別の細胞構成要素モジュールは、細胞構成要素モジュールAが含有しない他の細胞構成要素を更に含有することを条件として、細胞構成要素1、3及び10を含有し得る。上記の例示的な細胞構成要素モジュールAを更に考慮すると、細胞構成要素モジュールのセットにおける別の細胞構成要素モジュールは、細胞構成要素モジュールAが含有しない他の細胞構成要素を更に含有するという要件なしに提供される細胞構成要素1、3及び10のサブセットに限定され得る（しかしながら、そのような追加の細胞構成要素も有し得る）。

40

【0258】

いくつかの実施形態において、細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々は、複数の細胞構成要素のそれぞれの独立したサブセットにおける同じ

50

又は異なる数の細胞構成要素を含む。いくつかの実施形態において、それぞれの細胞構成要素モジュールの各々に対応する細胞構成要素のそれぞれの独立したサブセットの各々は、細胞構成要素の固有のサブセットである（例えば、非重複であり、複数の細胞構成要素における細胞構成要素の各々は、1つ以下のモジュールにグループ化される）。いくつかの実施形態において、第1の細胞構成要素モジュールは、第2の細胞構成要素モジュールに対応する細胞構成要素の第2のサブセットと重複する細胞構成要素の第1のサブセットを有する（例えば、重複しており、複数の細胞構成要素における少なくとも1つの細胞構成要素は、2つ以上の異なるモジュールに共通している）。

#### 【0259】

図3Cのブロック330を参照すると、いくつかの実施形態において、それぞれの細胞構成要素モジュールにおける複数の細胞構成要素の独立したサブセットは、5つ以上の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールにおける複数の細胞構成要素の独立したサブセットは、少なくとも2個、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも25個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも600個、少なくとも700個、少なくとも800個、少なくとも900個、少なくとも1000個、少なくとも2000個、又は少なくとも3000個の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素の独立したサブセットは、5000個以下、3000個以下、1000個以下、500個以下、200個以下、100個以下、90個以下、80個以下、70個以下、60個以下、又は50個以下の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素の独立したサブセットは、5～100個、2～300個、20～500個、200～1000個、又は1000～5000個の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素の独立したサブセットは、2個以上の細胞構成要素から始まり、5000個以下の細胞構成要素で終わる別の範囲内にある。

#### 【0260】

いくつかの実施形態において、それぞれの細胞構成要素モジュールにおける複数の細胞構成要素の独立したサブセットは、目的の生理学的状態と関連付けられた細胞プロセス（例えば、分子経路）における細胞構成要素からなる。例えば、図3Cのブロック332を参照すると、いくつかの実施形態において、それぞれの細胞構成要素モジュールにおける複数の細胞構成要素の独立したサブセットは、目的の生理学的状態と関連付けられた分子経路における2～20個の細胞構成要素からなる。

#### 【0261】

図3Dのブロック334を参照すると、細胞構成要素モジュールのセットにおける少なくとも第1の細胞構成要素モジュールは、目的の生理学的状態と関連付けられる。実際に、多数の細胞構成要素モジュールは、目的の生理学的状態と関連付けられ得る。

#### 【0262】

図3Dのブロック336を参照すると、1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、細胞構成要素モジュールのセットにおける対応する細胞構成要素モジュールを表す。

#### 【0263】

図3Dのブロック338を参照すると、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々は、複数の細胞構成要素の独立したサブセットを含む。

#### 【0264】

図3Dのブロック340を参照すると、いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールである。第1の細胞構成要素モジュールを含む複数の細胞構成要素モジュールの第1のサブセットは、目的の生理学的状態と

10

20

30

40

50

関連付けられる。すなわち、このような細胞構成要素モジュールは、目的の生理学的状態に  
関与する細胞構成要素を表す。例えば、このような細胞構成要素モジュールのこのよう  
な細胞構成要素は、いくつかのベースライン、野生型状態の細胞と比較して、目的の生理  
学的状態を表す細胞において下方制御又は上方制御され得る。更に、複数の細胞構成要素  
モジュールの第2のサブセットは、目的の生理学的状態と関連付けられていない。すなわ  
ち、このような細胞構成要素モジュールの細胞構成要素は、目的の生理学的状態に関与し  
ていない細胞構成要素を表す。例えば、このような細胞構成要素は、いくつかのベースラ  
イン、野生型状態の細胞と比較して、目的の生理学的状態を表す細胞において下方制御又  
は上方制御されない。このような実施形態において、（細胞構成要素モジュールの第1の  
サブセット内にある）第1の細胞構成要素モジュールについてのそれぞれの計算された活  
性化スコアが第1の閾値基準を満たし、複数の細胞構成要素モジュールの第2のサブセッ  
トにおける細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが第2  
の閾値基準を満たす場合、化学化合物は、目的の生理学的状態と識別される。例示的な第  
1の閾値基準は、図3Eのブロック348に関して以下で説明される。一般に、求められるのは、  
（第1の閾値を満たす計算された活性化スコアを有することによって示されるように）細胞  
構成要素モジュールの第1のサブセットにおける細胞構成要素モジュールと識別するが、  
（第2の閾値を満たす計算された活性化スコアを有することによって示されるように）細胞  
構成要素モジュールの第2のサブセットにおける細胞構成要素モジュールと識別しない  
化学化合物である。例えば、いくつかの実施形態において、第1の閾値の達成は、第1の  
所定の数値を上回る活性化スコアを必要とするが、第2の閾値の達成は、第2の所定  
の数値を下回る活性化スコアを必要とし、正確な第1及び第2の所定の数値は、適用に  
依存する。

#### 【0265】

上記に示されるように、いくつかの実施態様において、方法は、1つ以上のタイプの分子  
データ（例えば、細胞構成要素）を使用して、目的の生理学的状態（例えば、細胞プロ  
セス）を特徴付けることを含む。そのような分子データは、オミクスプロファイリング（  
例えば、トランスクリプトミクス、プロテオミクス、メタボロミクスなど）などの測定可  
能な属性（例えば、存在量及び/又は発現レベル）を有する任意の分析物を含むことが  
できる。

#### 【0266】

一般に、細胞プロセスと関連付けられる場合、細胞構成要素（例えば、遺伝子）の細胞  
構成要素モジュールは、同様の時間にスイッチする細胞構成要素（例えば、遺伝子）が一  
緒にモジュールを形成する、一連のスイッチングイベントから生じると考えられ得る。し  
たがって、例えば、いくつかの実施形態において、それぞれの細胞構成要素モジュールは  
、複数の細胞構成要素のそれぞれのサブセットを含み、細胞構成要素のサブセットは、目  
的のそれぞれの生理学的状態（例えば、目的の細胞プロセス）と関連付けられた挙動の類  
似性に基づいてグループ化される。一例では、目的のそれぞれの生理学的状態と関連付  
けられた細胞構成要素モジュールは、それぞれの生理学的状態を有する複数の細胞型にわた  
って同様に挙動する（例えば、同様の発現プロファイルを示す）遺伝子のサブセットを含  
むことができる。

#### 【0267】

図3Dのブロック342を参照すると、いくつかの実施形態において、細胞構成要素モ  
ジュールのセットは、第1の細胞構成要素モジュールからなる。

#### 【0268】

図3Dのブロック344を参照すると、いくつかの実施形態において、細胞構成要素モ  
ジュールのセットは、5つ以上の細胞構成要素モジュールを含む。いくつかの実施形態に  
おいて、細胞構成要素モジュールのセットは、20個以上、30個以上、40個以上、5  
0個以上、60個以上、70個以上、80個以上、90個以上、又は100個以上の細胞  
構成要素モジュールを含む。

#### 【0269】

10

20

30

40

50

いくつかの実施形態において、細胞構成要素モジュールのセットは、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも25個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも600個、少なくとも700個、少なくとも800個、少なくとも900個、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、又は少なくとも5000個の細胞構成要素モジュールを含む。いくつかの実施形態において、細胞構成要素モジュールのセットは、10,000個以下、5000個以下、2000個以下、1000個以下、500個以下、300個以下、200個以下、100個以下、90個以下、80個以下、70個以下、60個以下、又は50個以下の細胞構成要素モジュールを含む。いくつかの実施形態において、細胞構成要素モジュールのセットは、10~2000個の細胞構成要素モジュールからなる。いくつかの実施形態において、細胞構成要素モジュールのセットは、50~500個の細胞構成要素モジュールからなる。いくつかの実施形態において、細胞構成要素モジュールのセットは、5~20個、20~50個、50~100個、100~200個、200~500個、500~1000個、1000~5000個、又は5000~10,000個の細胞構成要素モジュールを含む。いくつかの実施形態において、細胞構成要素モジュールのセットは、5個以上の細胞構成要素モジュールから始まり、10,000個以下の細胞構成要素モジュールで終わる別の範囲内にある。

10

20

【0270】

いくつかの実施形態において、方法は、目的の生理学的状態と関連付けられた細胞構成要素モジュールを識別することを更に含む。そのような方法は、図14A~図14Dと併せて、細胞構成要素モジュールの識別と題されたセクションで以下に説明される。

【0271】

活性化スコア。

図3Bのブロック308に記載されるように、モデルは、モデルへのフィンガープリントの入力にตอบสนองして、1つ以上の計算された活性化スコアを出力する。一般に、訓練されたモデル(ブロック308のモデル)の出力は、標識(例えば、数値的活性化スコア)を含む訓練データセット上で学習するプロセスを通じて定義され、訓練されたモデルの出力が検証ステップなどを介して性能の最小レベルを満たすまで、複数のパラメータを調整する。訓練モデルは、「モデル訓練」と題されたセクションで以下に更に開示される。

30

【0272】

いくつかの実施形態において、1つ以上の計算された活性化スコアにおける活性化スコアは、それぞれの化合物に対応するそれぞれの細胞構成要素モジュールについてのそれぞれの活性化重みである。例えば、いくつかの実施形態において、活性化スコアは、図2A~図2B及び図14A~図14Dを参照して、並びに図5の活性化データ構造に示される「細胞構成要素モジュールの識別」と題された以下のセクションに記載されるように得られる活性化重みであり、活性化スコアは、それぞれの(例えば、第1の)細胞構成要素モジュールの活性化(例えば、誘導及び/又は差次的発現)を示し、それぞれの化合物による治療に相関する及び/又は応答する。

40

【0273】

したがって、いくつかのそのような実施形態において、訓練されたモデルは、出力として、試験化学化合物と目的の生理学的状態(例えば、目的の生理学的状態と関連付けられた第1の細胞構成要素モジュール)との関連性を示す計算された活性化スコアを提供する。次いで、図3Eのブロック348を参照すると、方法は、第1の細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態で識別する(例えば、関連性を決定する)ことを含む。

【0274】

図3Eのブロック350を参照すると、いくつかの実施形態において、第1の閾値基準は、第1の細胞構成要素モジュールが閾値活性化スコアを有することが必要である。一般

50



的に、求められるのは、（第1の閾値を満たす計算された活性化スコアを有することによって示されるように）目的の生理学的状態で識別する化学化合物である。例えば、いくつかの実施形態において、第1の閾値の達成は、第1の所定の数値を超える活性化スコアを必要とする。

**【0275】**

例えば、いくつかの実施形態において、活性化スコアは、「0」と「1」との間の正規化された連続値（又はA及びBが2つの異なる数である場合、いくつかの他の範囲の「A」から「B」）として表され、ここで、「1」に近い値（例えば、0.89、0.90、0.91、0.92など）は、細胞構成要素モジュール（及び細胞構成要素モジュールが表す化学化合物）と目的の生理学的状態との間の強い関連性を示す。「0」に近い値（例えば、0.01、0.02、0.03、0.04など）は、細胞構成要素モジュール（及び細胞構成要素が表す化学化合物）と目的の生理学的状態との間に関連性がないことを示す。そのような例では、第1の閾値は、「0」と「1」（又はA及びBが2つの異なる数であるいくつかの他の範囲の「A」から「B」）との間で選択され、細胞構成要素モジュール（及びそれが表す化学構造）は、活性化スコアが第1の閾値を上回る場合に目的の生理学的状態と関連付けられているとみなされ、一方、細胞構成要素モジュール（及びそれが表す化学構造）は、活性化スコアが第1の閾値を下回る場合に目的の生理学的状態と関連付けられていないとみなされる。いくつかのそのような実施形態において、活性化スコアは、「0」と「1」（又はA及びBが2つの異なる数であるいくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、第1の閾値は、0と1との間、0.10と0.90との間、0.20と0.80との間、0.30と0.70との間、0.50と0.99との間、0.60と0.99との間、0.70と0.99との間、0.80と0.99との間、又は0.90と0.99との間の値である。

10

20

**【0276】**

別の例として、いくつかの実施形態において、活性化スコアは、「0」と「1」（又はA及びBが2つの異なる数である場合、いくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、「1」に近い値（例えば、0.89、0.90、0.91、0.92など）は、細胞構成要素モジュール（及び細胞構成要素モジュールが表す化学化合物）と目的の生理学的状態との間に関連性がないことを示す。「0」に近い値（例えば、0.01、0.02、0.03、0.04など）は、細胞構成要素モジュール（及び細胞構成要素が表す化学化合物）と目的の生理学的状態との間の関連性を示す。そのような例では、第1の閾値は、「0」と「1」（又はA及びBが2つの異なる数であるいくつかの他の範囲の「A」から「B」）との間で選択され、細胞構成要素モジュール（及びそれが表す化学構造）は、活性化スコアが第1の閾値を下回る場合に目的の生理学的状態と関連付けられているとみなされ、一方、細胞構成要素モジュール（及びそれが表す化学構造）は、活性化スコアが第1の閾値を上回る場合に目的の生理学的状態と関連付けられていないとみなされる。いくつかのそのような実施形態において、活性化スコアは、「0」と「1」（又はA及びBが2つの異なる数であるいくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、第1の閾値は、0と1との間、0.10と0.90との間、0.20と0.80との間、0.30と0.70との間、0.50と0.99との間、0.60と0.99との間、0.70と0.99との間、0.80と0.99との間、又は0.90と0.99との間の値である。

30

40

**【0277】**

図3Eのブロック352を参照すると、いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュール（例えば、2～1000個、10～100個、2～100個、4～50個の細胞構成要素モジュール）であり、ブロック348の識別は、細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの計算された活性化スコアが第1の閾値基準を満たすことを必要とする。例えば、細胞構成要素モジュールのセットが2つの細胞構成要素モジュール：A及びBからなる

50

場合を考慮する。図 3 E のブロック 3 5 2 は、細胞構成要素モジュール A 及び B の活性化スコアの各々が、第 1 の閾値条件を満たすことを必要とする。例えば、細胞構成要素モジュール A が 0 . 2 5 の計算された活性化スコアを有し、細胞構成要素モジュール B が 0 . 7 5 の計算された活性化スコアを有し、第 1 の閾値条件の達成が、各活性化スコアが 0 . 4 よりも大きいことを必要とする場合を考慮する。この例では、各活性化スコアが 0 . 4 の閾値要件を超えないため、細胞構成要素モジュールのセットは、図 3 E のブロック 3 5 2 の要件を満たさない。

【 0 2 7 8 】

図 3 E のブロック 3 5 4 を参照すると、いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュール（例えば、2 ~ 1 0 0 0 個、1 0 ~ 1 0 0 個、2 ~ 1 0 0 個、4 ~ 5 0 個の細胞構成要素モジュール）であり、ブロック 3 4 8 の識別は、細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの計算された活性化スコアにわたる中心傾向の測定値が第 1 の閾値基準を満たすことを必要とする。例えば、細胞構成要素モジュールのセットが 2 つの細胞構成要素モジュール：A 及び B からなる場合を考慮する。図 3 E のブロック 3 5 4 は、細胞構成要素モジュール A 及び B の活性化スコアの中心傾向のいくつかの測定値が第 1 の閾値条件を満たすことを必要とする。例えば、中心傾向の測定値が平均化されており、細胞構成要素モジュール A が 0 . 2 5 の計算された活性化スコアを有し、細胞構成要素モジュール B が 0 . 7 5 の計算された活性化スコアを有し、第 1 の閾値条件の達成が、平均活性化スコアが 0 . 4 よりも大きいことを必要とする場合を考慮する。この例では、細胞構成要素モジュールのセットは、それらが 0 . 4 の閾値要件よりも大きい  $0 . 2 5 + 0 . 7 5 / 2$  又は  $0 . 5$  の平均活性化スコアを有するため、図 3 E のブロック 3 5 4 の要件を満たす。いくつかの実施形態において、中心傾向の測定値は、細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの計算された活性化スコアの各々の算術平均、加重平均、ミッドレンジ、ミッドヒンジ、トリミアン、ウィンザライズド平均、平均、又はモードである。

【 0 2 7 9 】

化合物。

いくつかの実施形態において、試験化学化合物は、小分子、生物製剤、タンパク質、小分子と組み合わせられたタンパク質、ADC、siRNA 若しくは干渉RNAなどの核酸、cDNA 過剰発現野生型及び/若しくは変異体shRNA、cDNA 過剰発現野生型及び/若しくは変異体ガイドRNA（例えば、Cas9系若しくは他の細胞成分編集系）、並びに/又は前述のいずれかの任意の組み合わせである。

【 0 2 8 0 】

いくつかの実施形態において、試験化学化合物は、無機又は有機である。

【 0 2 8 1 】

例えば、図 3 E のブロック 3 5 6 を参照すると、いくつかの実施形態において、試験化学化合物は、2 0 0 0 ダルトン (Da) 未満の分子量を有する有機化合物である。いくつかの実施形態において、試験化学化合物は、少なくとも 1 0 Da、少なくとも 2 0 Da、少なくとも 5 0 Da、少なくとも 1 0 0 Da、少なくとも 2 0 0 Da、少なくとも 5 0 0 Da、少なくとも 1 k Da、少なくとも 2 k Da、少なくとも 3 k Da、少なくとも 5 k Da、少なくとも 1 0 k Da、少なくとも 2 0 k Da、少なくとも 3 0 k Da、少なくとも 5 0 k Da、少なくとも 1 0 0 k Da、又は少なくとも 5 0 0 k Da の分子量を有する。いくつかの実施形態において、試験化学化合物は、1 0 0 0 k Da 以下、5 0 0 k Da 以下、1 0 0 k Da 以下、5 0 k Da 以下、1 0 k Da 以下、5 k Da 以下、2 k Da 以下、1 k Da 以下、5 0 0 Da 以下、3 0 0 Da 以下、1 0 0 Da 以下、又は 5 0 Da 以下の分子量を有する。いくつかの実施形態において、試験化学化合物は、1 0 Da ~ 9 0 0 Da、5 0 Da ~ 1 0 0 0 Da、1 0 0 Da ~ 2 0 0 0 Da、1 k Da ~ 1 0 k Da、5 k Da ~ 5 0 0 k Da、又は 1 0 0 k Da ~ 1 0 0 0 k Da の分子量を有する。いくつかの実施形態において、試験化学化合物は、1 0 ダルトン以上から始まり、1 0 0 0 k D

a 以下で終わる別の範囲内にある分子量を有する。

【0282】

図3Eのブロック358を参照すると、いくつかの実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則の各々を満たす有機化合物である。5つ（例えば、RO5）の基準のリピンスキーの法則は、それぞれの薬理的又は生物学的活性を有するそれぞれの化合物が、ヒトへの投与に好適な対応する化学的又は物理的特性を有するかどうかを決定するなどの、ドラッグライクネスを評価するために使用されるガイドラインのセットである。5つのリピンスキーの法則は、化合物のドラッグライクネスを決定するための以下の基準を含む。(i) 500Da未満の分子量、(ii) 5個以下の水素結合ドナー、(iii) 10個以下の水素結合アクセプター、及び(iv) 5個以下のオクタノール-水分配係数  $\log P$ 。

【0283】

図3Eのブロック360を参照すると、いくつかの実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則の少なくとも2つ、3つ、又は4つの基準を満たす有機化合物である。いくつかの実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則のゼロ、1つ、2つ、3つ、又は4つ全ての基準を満たす有機化合物である。

【0284】

いくつかの実施形態において、試験化学化合物は、データベースから選択される。薬物スクリーニング、注釈、及び/又は化合物標的及び化合物の化学特性などの一般的な情報からの結果を提供する好適な化合物データベースの例としては、限定されないが、Genomics of Drug Sensitivity in Cancer、Cancer Therapeutics Response Portal、Connectivity Map、PharmacDB、Base of Bioisosterically Exchangeable Replacements (BoBER)、及び/又はDrugBankが挙げられる。いくつかの実施形態において、試験化学化合物は、遺伝子及び遺伝子産物、摂動誘発細胞構成要素シグネチャ、及び/又は経路注釈に関する情報を提供するデータベースから選択される。好適なデータベースの例としては、限定されないが、NIH Gene Expression Omnibus (GEO)、EBI ArrayExpress、NCBI、BLAST、EMBL-EBI、GenBank、Ensembl、KEGG経路データベース、Library of Integrated Network-based Cellular Signatures (LINCS) L1000データセット、Reactome経路データベース、及び/又はGene Ontologyプロジェクトが挙げられる。

【0285】

方法300の結果を実際の適用に使用する。

いくつかの実施形態において、図3と併せて上述した方法300を使用して、目的の生理学的状態に対して複数の試験化合物を評価する。そのような実施形態において、複数の試験化合物における試験化合物の各々は、図3の方法300により実行される。したがって、100個の試験化合物及び1つの目的の生理学的状態が存在する場合、そのような実施形態において、方法300を100回実行し、100回の各事例は、試験化合物の異なる1つに対してである。

【0286】

更に、いくつかの実施形態において、図3と併せて上述した方法300を使用して、目的の複数の生理学的状態に対して複数の化合物を評価する。そのような実施形態において、目的の生理学的状態の各々について、複数の試験化合物におけるそれぞれの各試験化合物の各々は、図3の方法300により実行される。したがって、100個の試験化合物及び2つの目的の生理学的状態が存在する場合、そのような実施形態において、方法300は200回実行され、200回の各事例は、目的の第1の生理学的状態又は第2の生理学的状態のいずれかに対する試験化合物の異なる1つに対してである。

10

20

30

40

50

## 【0287】

いくつかの実施形態において、複数の試験化合物は、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも800個、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、少なくとも5000個、少なくとも8000個、少なくとも10,000個、少なくとも20,000個、少なくとも30,000個、少なくとも50,000個、少なくとも80,000個、少なくとも100,000個、少なくとも200,000個、少なくとも500,000個、少なくとも800,000個、少なくとも100万個、又は少なくとも200万個の試験化合物を含み、単一の目的の生理学的状態が存在する。いくつかのそのような実施形態において、方法300は、少なくとも5回、少なくとも10回、少なくとも15回、少なくとも20回、少なくとも30回、少なくとも40回、少なくとも50回、少なくとも100回、少なくとも200回、少なくとも300回、少なくとも400回、少なくとも500回、少なくとも800回、少なくとも1000回、少なくとも2000回、少なくとも3000回、少なくとも4000回、少なくとも5000回、少なくとも8000回、少なくとも10,000回、少なくとも20,000回、少なくとも30,000回、少なくとも50,000回、少なくとも80,000回、少なくとも100,000回、少なくとも200,000回、少なくとも500,000回、少なくとも800,000回、又は少なくとも200万回実行されて、少なくとも5、少なくとも10、少なくとも15、少なくとも20、少なくとも30、少なくとも40、少なくとも50、少なくとも100、少なくとも200、少なくとも300、少なくとも400、少なくとも500、少なくとも800、少なくとも1000、少なくとも2000、少なくとも3000、少なくとも4000、少なくとも5000、少なくとも8000、少なくとも10,000、少なくとも20,000、少なくとも30,000、少なくとも50,000、少なくとも80,000、少なくとも100,000、少なくとも200,000、少なくとも500,000、少なくとも800,000、又は少なくとも100万、又は少なくとも200万の活性化スコアを実現し、各試験化合物に対して1つである。

## 【0288】

いくつかの実施形態において、複数の化合物は、1000万個以下、500万個以下、100万個以下、500,000個以下、100,000個以下、50,000個以下、10,000個以下、8000個以下、5000個以下、2000個以下、1000個以下、800個以下、500個以下、200個以下、又は100個以下の試験化合物を含む。いくつかの実施形態において、複数の化合物は、10~500個、100~10,000個、5000~200,000個、又は10,000~100万個の試験化合物からなる。

## 【0289】

いくつかの実施形態において、複数の試験化合物は、10~ $1 \times 10^6$ 個の試験化合物である。いくつかの実施形態において、複数の試験化合物は、100~100,000個の試験化合物である。いくつかの実施形態において、複数の試験化合物は、1000~100,000個の試験化合物である。

## 【0290】

したがって、方法300を使用して、多数の試験化合物についての活性化スコアを得ることができる。これらの活性化スコアに対する第1の閾値の適用を使用して、目的の生理学的状態と関連付けられる試験された多くの試験化合物の中から、試験化合物を識別することができる。典型的な実施形態において、選択された数の試験化合物は、それらが目的の生理学的状態と関連付けられることを示す活性化スコアを有するが、他のものはそうではない。選択された数の試験化合物の分析を使用して、目的の生理学的状態との関連性をもたらす化合物を試験するための分子特性を決定することができる。例えば、目的の生理学的状態と関連付けられていることを示す活性化スコアを有する選択された数の試験化合

物の化学構造を、目的の生理学的状態と関連付けられていない試験化合物と区別する構造の類似性について視覚的に検査することができる。次いで、そのような分子特性は、モデル601によって評価された元の試験分子に含まれず、モデル601を訓練するために使用されなかった新しい試験分子に組み込むことができる。

**【0291】**

更に、より正式なアプローチを使用して、試験化合物を分析することができる（方法300によって課された第1の閾値を満たすものと満たさないものとの両方）。例えば、部分構造マイニングを使用して、そのような化合物を目的の生理学的状態と関連付けるようにする試験化合物内の部分構造を特定することができる。部分構造マイニングの例としては、MOSS（参照により本明細書に組み込まれる、Borgetl and Meinl, 2006, "Full Perfect Extension Pruning for Frequent Graph Mining," Proc. Workshop on Mining Complex Data (MCD 2006 at ICDM 2006, Hong Kong, China, IEEE Press, Piscataway, NJ, USA、及びMOFA（参照により本明細書に組み込まれる、Meinl and Worlein, 2006 "Mining Molecular Datasets on Symmetric Processor Systems," International conference on Systems, man and Cybernetics 2, pp. 1269 - 1274）が挙げられるが、これらに限定されない。

10

**【0292】**

また、最大共通部分構造（MCS）分析を使用して、そのような化合物を目的の生理学的状態と関連付ける試験化合物内の部分構造を識別することができる。MCS分析の例としては、LIBMCS（Chemaxon, Library MCS, 2008）、MCSS（OEChem TK version 2.0.0, OpenEye Scientific Software, Santa Fe, NM. <http://www.eyesopen.com>）、及びCncMCS（<http://www.chemnavigator.com/cnc/products/downloads.asp>）が挙げられるが、これらに限定されない。

20

**【0293】**

また、SMARTSを使用して、そのような化合物を目的の生理学的状態と関連付ける試験化合物内の部分構造を識別することができる。SMART分析の例は、CDK Descriptor GUIである。

30

**【0294】**

また、頻出部分グラフマイニングを使用して、そのような化合物を目的の生理学的状態と関連付けるようにする試験化合物内の部分構造を識別することができる。頻出部分グラフマイニングの例は、ParMol（Uni Erlangen）である。

**【0295】**

また、グラフ及び化学マイニングを使用して、そのような化合物を目的の生理学的状態と関連付けるようにする試験化合物内の部分構造を識別することができる。グラフ及び化学マイニングの例は、PAFI/AFGen（Karypis Lab UMN）である。

40

**【0296】**

摂動シグネチャ。

上記のように、いくつかの実施形態において、目的の生理学的状態は、摂動シグネチャ（例えば、摂動に回答して第1の細胞状態と第2の細胞状態との間の不一致を特徴とする）である。したがって、本開示の別の態様は、試験化学化合物を目的の生理学的状態と関連付ける方法700を提供する。いくつかの実施形態において、目的の生理学的状態は、疾患である。

**【0297】**

ブロック702を参照すると、方法は、試験化学化合物の化学構造のフィンガープリントを得ることを含む。「生理学的状態」及び「化合物」と題する上記のセクションに開示

50

されるような、生理学的状態、化合物、フィンガープリント、及び/又はフィンガープリントを得る方法の任意の好適な実施形態は、当業者には明白であろうように、それらの任意の置換、修飾、追加、欠失、及び/又は組み合わせを含むことが企図される。

**【0298】**

例えば、いくつかの実施形態において、試験化学化合物は、2000ダルトン未満の分子量を有する有機化合物である。いくつかの実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則の各々を満たす有機化合物である。いくつかの実施形態において、試験化学化合物は、5つの基準のリピンスキーの法則のうち少なくとも3つの基準を満たす有機化合物である。いくつかの実施形態において、方法は、試験化学化合物の単純化された分子入力ラインエントリーシステム(SMILES)文字列表現からフィンガープリントを計算することを更に含む。いくつかの実施形態において、フィンガープリントは、SMILES Transformer、ECFP4、RNNS2S、又はGraphConvを使用して、化学構造から生成される。

10

**【0299】**

ブロック704を参照すると、方法は、フィンガープリントをモデルに入力することを更に含み、モデルは100以上のパラメータを含み、モデルは、フィンガープリントのモデルへの入力に回答して1つ以上の計算された活性化スコアを出力し、1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、摂動シグネチャのセットにおける対応する摂動シグネチャを表す。

**【0300】**

「モデルアーキテクチャ」と題された上記のセクションに開示されるものなどのモデルの任意の好適な実施形態が企図され、当業者には明らかであろうように、それらの任意の置換、修飾、追加、削除、及び/又は組み合わせが企図される。例えば、いくつかの実施形態において、モデルは、ニューラルネットワークを含む。いくつかのそのような実施形態において、ニューラルネットワークは、ReLU活性化を有する完全に接続されたニューラルネットワークである。いくつかの実施形態において、ニューラルネットワークは、メッセージパッシングニューラルネットワークである。

20

**【0301】**

いくつかの実施形態において、モデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

30

**【0302】**

いくつかの実施形態において、モデルは、複数のコンポーネントモデルのアンサンブルモデルであり、1つ以上の計算された活性化スコアにおける計算された活性化スコアの各々は、複数のコンポーネントモデルにおけるコンポーネントモデルの各々の出力の中心傾向の測定値である。

**【0303】**

いくつかの実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

40

**【0304】**

いくつかの実施形態において、複数のコンポーネントモデルは、複数のニューラルネットワークを含む。いくつかのそのような実施形態において、複数のニューラルネットワークにおける第1のニューラルネットワークは、ReLU活性化を伴う完全に接続されたニューラルネットワークであり、複数のニューラルネットワークにおける第2のニューラルネットワークは、メッセージパッシングニューラルネットワークである。

**【0305】**

上記で定義されるように、摂動は、1つ以上の化合物による治療などの1つ以上の状態

50

への細胞の任意の曝露を指す。いくつかの実施形態において、摂動シグネチャは、摂動によって誘発される細胞内の1つ以上の細胞構成要素の発現又は存在量レベルの変化である。

#### 【0306】

例示的な摂動には、限定されないが、遺伝子ノックダウン、刺激に対する細胞応答、組織成長及び再生、並びに/又は化合物による治療若しくは化合物への曝露が含まれる。例示的なペルターバゲンには、小分子、生物製剤、治療剤、タンパク質、小分子と組み合わされたタンパク質、ADC、siRNA若しくは干渉RNAなどの核酸、cDNA過剰発現野生型及び/若しくは変異体shRNA、cDNA過剰発現野生型及び/若しくは変異体ガイドRNA（例えば、Cas9系若しくは他の遺伝子編集系）、又は前述のいずれかの任意の組み合わせが含まれるが、これらに限定されない。

10

#### 【0307】

いくつかの実施形態において、摂動は、システムレベル（例えば、結合又はドッキング活性）、並びに/又は下流効果及び臓器レベルの表現型に関して特徴付けられる。いくつかの実施形態において、摂動は、分子、細胞、及び/又は組織レベルでのペルターバゲンに対する応答を駆動する又はその基礎となる機構の機能として特徴付けられる（例えば、摂動の前又は後にバイオマーカー、細胞生存率、及び/又は薬物タンパク質相互作用を識別又は測定することによって）。例えば、摂動の測定値は、表現型の測定値（例えば、IC50値）及び/又は細胞構成要素シグネチャ（例えば、オミクスプロファイリング）を含むことができる。

20

#### 【0308】

いくつかの実施形態において、それぞれの摂動及び/又は対応する摂動シグネチャは、Genomics of Drug Sensitivity in Cancer、Cancer Therapeutics Response Portal、Connectivity Map、PharmacoDB、Base of Bioisosterically Exchangeable Replacements (BOBER)、DrugBank、Human Cell Atlas、Molecular Signatures Database (MSigDB)、及び/又はEnrichrなどの公的に利用可能なデータベースから得られる。摂動データを得ることができる他の好適なデータベースとしては、NIH Gene Expression Omnibus (GEO)、EBI ArrayExpress、NCBI、BLAST、EMBL-EBI、GenBank、Ensembl、KEGG経路データベース、Library of Integrated Network-based Cellular Signatures (LINCS) L1000データセット、Reactome経路データベース、及び/又はGene Ontologyプロジェクトが挙げられる。

30

#### 【0309】

摂動データを得る方法には、例えば、perturb-seq、CRISP-seq、CROP-seq、CRISPRi、TAP-seq、CRISPRa、perturb-CITE-seq、sci-Plex、multiplexed、MIX-seq、CyTOF、及び/又はscRNA-seqを使用した細胞構成要素データの測定が含まれる。摂動データを得る方法には、更に、質量分析（例えば、LCMS、GCMS）、フローサイトメトリー、定量的ポリメラーゼ連鎖反応（qPCR）、ゲル電気泳動、遺伝子チップ分析、マイクロアレイ、細胞蛍光分析、蛍光顕微鏡、共焦点レーザーキャニング顕微鏡、レーザーキャニングサイトメトリー、親和性クロマトグラフィー、手動バッチモード分離、電界懸濁、配列決定、及び/又はそれらの任意の組み合わせを含む、オミクスデータを得る任意の方法が含まれる。いくつかの実施形態において、本明細書に開示される細胞構成要素存在量値を得るための方法のうちいずれかは、摂動データを得る際に（例えば、摂動シグネチャのために）使用するために企図される。

40

#### 【0310】

いくつかの実施形態において、摂動シグネチャのセットは、第1の摂動シグネチャから

50

なる。いくつかの実施形態において、摂動シグネチャのセットは、5つ以上の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、10個以上の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、100個以上の摂動シグネチャを含む。

【0311】

いくつかの実施形態において、摂動シグネチャのセットは、少なくとも2つ、少なくとも3つ、少なくとも4つ、少なくとも5つ、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも800個、少なくとも1000個、少なくとも2000個、又は少なくとも5000個の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、10,000個以下、5000個以下、1000個以下、800個以下、500個以下、200個以下、100個以下、50個以下、又は20個以下の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、5~50個、2~100個、20~500個、10~1000個、800~5000個、又は50~2000個の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、2つ以上の摂動シグネチャから始まり、10,000個以下の摂動シグネチャで終わる別の範囲内にある。

10

【0312】

ブロック706を参照すると、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。

20

【0313】

いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動状態は、複数の化合物における化合物に曝露されていない対照細胞によって表される。いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動状態は、それぞれの摂動シグネチャに関連付けられた化合物以外の複数の化学化合物における化学化合物に曝露されている無関係の摂動細胞にわたる平均によって表される。

30

【0314】

いくつかの実施形態において、細胞状態の変化は、変化していない細胞状態と変化した細胞状態との間の変化を指し、変化した細胞状態は、変化していない細胞状態から変化した細胞状態への細胞遷移を通じて生じる。更に、(i)変化していない細胞状態、(ii)変化した細胞状態、及び(iii)変化していない細胞状態から変化した細胞状態への遷移のうち少なくとも1つが、目的の生理学的状態と関連付けられる。

40

【0315】

いくつかの実施形態において、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャは、非限定的な例として、参照により本明細書に組み込まれる、2019年7月15日に出願された「Methods of Analyzing Cells」と題された米国特許出願第16/511,691号に開示された方法のいずれかを使用して決定され得る。

【0316】

ある特定の実施形態において、摂動(例えば、特定の化学組成物への細胞の曝露)の共変量が存在し得る。例えば、化学組成物の共変量は、化学組成物の特定の用量、化学組成

50



物に曝露された細胞が細胞構成要素を定量化するために測定される時間、及び/又は化学組成物に曝露された細胞の同一性（例えば、細胞株）を含み得る。いくつかの実施形態において、摂動（例えば、特定の化学組成物への細胞の曝露）は、その共変量の閾値量も特定の細胞遷移に影響すると予測される場合にのみ、特定の細胞遷移に影響すると予測される。言い換えれば、いくつかの実施形態において、特定の摂動シグネチャの計算された活性化スコアは、特定の摂動シグネチャの化学組成物の共変量が、目的の生理学的状態と関連付けられた特定の細胞遷移にも影響を与えると予測されるかどうかによって少なくとも部分的に決定される。

#### 【0317】

一般に、上述のように、訓練されたモデルの出力は、標識（例えば、数値的活性化スコア）を含む訓練データセット上で学習するプロセスを通じて定義され、訓練されたモデルの出力が検証ステップなどを介して性能の最小レベルを満たすまで、複数のパラメータを調整する。訓練モデルは、「モデル訓練」と題されたセクションで以下に更に開示される。したがって、いくつかのそのような実施形態において、訓練されたモデルは、出力として、試験化学化合物と目的の生理学的状態との関連性を示す第1の摂動シグネチャについて計算された活性化スコアを提供する（例えば、第1の摂動シグネチャは、目的の生理学的状態と関連付けられた細胞状態遷移と関連付けられる）。

10

#### 【0318】

次にブロック708を参照すると、方法は、摂動シグネチャのセットにおける第1の摂動シグネチャについてのそれぞれの計算された活性化スコアが第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態と識別することを含む。

20

#### 【0319】

「活性化スコア」と題する上記のセクションに開示されるような活性化スコアの任意の好適な実施形態は、1つ以上の計算された活性化スコアを得るために企図され、活性化スコアの各々は、当業者に明らかであろうように、それらの任意の置換、修飾、追加、欠失、及び/又は組み合わせを含む、摂動シグネチャのセットにおける対応する摂動シグネチャを表す。

#### 【0320】

一般的に、求められるのは、（第1の閾値基準を満たす計算された活性化スコアを有することによって示されるように）目的の生理学的状態と識別する化学化合物である。例えば、いくつかの実施形態において、第1の閾値の達成は、第1の所定の数値を超える活性化スコアを必要とする。

30

#### 【0321】

例えば、いくつかの実施形態において、活性化スコアは、「0」と「1」（又はA及びBが2つの異なる数である場合、いくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、ここで、「1」に近い値（例えば、0.89、0.90、0.91、0.92など）は、摂動シグネチャ（及び摂動シグネチャが表す化学化合物）と目的の生理学的状態との間の強い関連性を示す。「0」に近い値（例えば、0.01、0.02、0.03、0.04など）は、摂動シグネチャ（及び摂動シグネチャが表す化学化合物）と目的の生理学的状態との間に関連性がないことを示す。そのような例では、第1の閾値は、「0」と「1」（又はA及びBが2つの異なる数であるいくつかの他の範囲の「A」から「B」）との間で選択され、摂動シグネチャ（及びそれが表す化学構造）は、活性化スコアが第1の閾値を上回る場合に目的の生理学的状態と関連付けられているとみなされ、一方、摂動シグネチャ（及びそれが表す化学構造）は、活性化スコアが第1の閾値を下回る場合に目的の生理学的状態と関連付けられていないとみなされる。いくつかのそのような実施形態において、活性化スコアは、「0」と「1」（又はA及びBが2つの異なる数であるいくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、第1の閾値は、0と1との間、0.10と0.90との間、0.20と0.80との間、0.30と0.70との間、0.50と0.99との間、0.60と0.99との間、0.70と0.99との間、0.80

40

50

と 0.99 との間、又は 0.90 と 0.99 との間の値である。

【0322】

別の例として、いくつかの実施形態において、活性化スコアは、「0」と「1」（又は A 及び B が 2 つの異なる数である場合、いくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、ここで、「1」に近い値（例えば、0.89、0.90、0.91、0.92 など）は、摂動シグネチャ（及び摂動シグネチャが表す化学化合物）と目的の生理学的状態との間に関連性がないことを示す。「0」に近い値（例えば、0.01、0.02、0.03、0.04 など）は、摂動シグネチャ（及び摂動シグネチャが表す化学化合物）と目的の生理学的状態との間の関連性を示す。そのような例では、第 1 の閾値は、「0」と「1」（又は A 及び B が 2 つの異なる数であるいくつかの他の範囲の「A」から「B」）との間で選択され、摂動シグネチャ（及びそれが表す化学構造）は、活性化スコアが第 1 の閾値を下回る場合に目的の生理学的状態と関連付けられているとみなされ、一方、摂動シグネチャ（及びそれが表す化学構造）は、活性化スコアが第 1 の閾値を上回る場合に目的の生理学的状態と関連付けられていないとみなされる。いくつかのそのような実施形態において、活性化スコアは、「0」と「1」（又は A 及び B が 2 つの異なる数であるいくつかの他の範囲の「A」から「B」）との間の連続的な尺度における正規化された値として表され、第 1 の閾値は、0 と 1 との間、0.10 と 0.90 との間、0.20 と 0.80 との間、0.30 と 0.70 との間、0.50 と 0.99 との間、0.60 と 0.99 との間、0.70 と 0.99 との間、0.80 と 0.99 との間、又は 0.90 と 0.99 との間の値である。

10

20

【0323】

いくつかの実施形態において、第 1 の閾値基準は、第 1 の摂動シグネチャが閾値活性化スコアを有することが必要である。

【0324】

いくつかの実施形態において、第 1 の閾値基準は、第 1 の摂動シグネチャが、摂動シグネチャのセットにおける少なくとも閾値ランクを有することが必要であり、摂動シグネチャのセットは、摂動シグネチャのセットにおける摂動シグネチャの各々と参照シグネチャ（例えば、単一細胞遷移シグネチャ）との比較に基づいてランク付けされる。化学化合物を生理学的状態と関連付ける際の使用に好適な参照シグネチャ（例えば、単一細胞遷移シグネチャ）に対する摂動シグネチャの比較方法は、以下の「摂動シグネチャについての数値的活性化スコア」と題されたセクションに更に詳細に記載されている。

30

【0325】

いくつかの実施形態において、識別は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアが閾値基準を満たすことを必要とする。いくつかの実施形態において、識別は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアにわたる中心傾向の測定値が閾値基準を満たすことを必要とする。いくつかの実施形態において、中心傾向の測定値は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアの各々の算術平均、加重平均、ミッドレンジ、ミッドヒンジ、トリミアン、ウィンザライズド平均、平均、又はモードである。

40

【0326】

いくつかの実施形態において、摂動シグネチャのセットは、2 ~ 100 個の摂動シグネチャであり、識別は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアが閾値基準を満たすことを必要とする。いくつかの実施形態において、摂動シグネチャのセットは、2 ~ 100 個の摂動シグネチャであり、識別は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアにわたる中心傾向の測定値が閾値基準を満たすことを必要とする。いくつかの実施形態において、中心傾向の測定値は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアの各々の算術平均、加重平均、ミッドレンジ、ミッドヒンジ、トリミアン、ウィンザライズド平均、平均、又はモードである。

50

## 【0327】

いくつかの実施形態において、摂動シグネチャのセットは、複数の摂動シグネチャであり、第1の摂動シグネチャを含む、複数の摂動シグネチャの第1のサブセットが、目的の生理学的状態と関連付けられ、複数の摂動シグネチャの第2のサブセットが、目的の生理学的状態と関連付けられておらず、第1の摂動シグネチャについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たし、複数の摂動シグネチャの第2のサブセットにおける摂動シグネチャについてのそれぞれの計算された活性化スコアが、第2の閾値基準を満たす場合、試験化学化合物が、目的の生理学的状態と識別される。

## 【0328】

いくつかの実施形態において、第2の閾値基準は、複数の摂動シグネチャの第2のサブセットにおける摂動シグネチャについてのそれぞれの計算された活性化スコアが閾値活性化スコアを有することを必要とする。 10

## 【0329】

いくつかの実施形態において、第2の閾値基準は、複数の摂動シグネチャの第2のサブセットにおける摂動シグネチャについてのそれぞれの計算された活性化スコアが、摂動シグネチャのセットにおける少なくとも閾値ランクを有することを必要とし、摂動シグネチャのセットは、摂動シグネチャのセットにおける摂動シグネチャの各々と参照シグネチャ（例えば、単一細胞遷移シグネチャ）との比較に基づいてランク付けされる。

## 【0330】

いくつかの実施形態において、識別は、摂動シグネチャの第2のサブセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアが第2の閾値基準を満たすことを必要とする。いくつかの実施形態において、識別は、摂動シグネチャの第2のサブセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアにわたる中心傾向の測定値が第2の閾値基準を満たすことを必要とする。いくつかの実施形態において、中心傾向の測定値は、摂動シグネチャのセットにおける摂動シグネチャの各々のそれぞれの計算された活性化スコアの各々の算術平均、加重平均、ミッドレンジ、ミッドヒンジ、トリミアン、ウィンザライズド平均、平均、又はモードである。 20

## 【0331】

III. 化学化合物を目的の生理学的状態と関連付ける方法  
モデル訓練。 30

本開示の別の態様は、化学化合物を目的の生理学的状態と関連付ける方法800を提供する。いくつかの実施形態において、目的の生理学的状態は、疾患である。

## 【0332】

ブロック802を参照すると、方法は、複数の化合物における化合物の各々の化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得る。「生理学的状態」及び「化合物」と題する上記のセクションに開示されるような、生理学的状態、化合物、フィンガープリント、及び/又はフィンガープリントを得る方法の任意の好適な実施形態は、当業者には明白であろうように、それらの任意の置換、修飾、追加、欠失、及び/又は組み合わせを含むことが企図される。

## 【0333】

例えば、いくつかの実施形態において、複数の化合物は、 $10 \sim 1 \times 10^6$ 個の化合物である。いくつかの実施形態において、複数の化合物は、 $100 \sim 100,000$ 個の化合物である。いくつかの実施形態において、複数の化合物は、 $1000 \sim 100,000$ 個の化合物である。 40

## 【0334】

いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、 $2000$ ダルトン未満の分子量を有する有機化合物である。いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則の各々を満たす。いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則のうち少なくとも3つの基準を満たす。いくつかの実施 50

形態において、それぞれのフィンガープリントの各々は、SMILES Transformer、ECFP4、RNNS2S、又はGraphConvを使用して、化学構造から生成される。

#### 【0335】

ブロック804を参照すると、方法は、複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることを含み、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む。上記「細胞構成要素及び細胞構成要素モジュール」、及び以下「細胞構成要素モジュールの識別」と題するセクションに開示されるように、細胞構成要素、細胞構成要素モジュール、及び/又は細胞構成要素モジュールを識別する方法の任意の好適な実施形態が企図され、当業者には明らかであるように、それらの任意の置換、修飾、追加、欠失、及び/又は組み合わせが含まれる。

10

#### 【0336】

例えば、いくつかの実施形態において、細胞構成要素モジュールのセットは、単一の細胞構成要素モジュールである。いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールである。いくつかの実施形態において、細胞構成要素モジュールのセットは、200~500個の細胞構成要素モジュールである。いくつかの実施形態において、細胞構成要素モジュールのセットは、単一の細胞構成要素モジュールからなる。いくつかの実施形態において、細胞構成要素モジュールのセットは、5つ以上の細胞構成要素モジュールを含む。いくつかの実施形態において、細胞構成要素モジュールのセットは、10個以上の細胞構成要素モジュールを含む。いくつかの実施形態において、細胞構成要素モジュールのセットは、100個以上の細胞構成要素モジュールを含む。いくつかの実施形態において、細胞構成要素モジュールのセットは、複数の細胞構成要素モジュールであり、複数の細胞構成要素モジュールの第1のサブセットは、目的の生理学的状態と関連付けられ、複数の細胞構成要素モジュールの第2のサブセットは、目的の生理学的状態と関連付けられていない。

20

#### 【0337】

いくつかの実施形態において、図2A~図2Bの例示的なワークフローによって示されるように、方法は、電子形式で1つ以上の第1のデータセットを得ることであって、1つ以上の第1のデータセットが、第1の複数の細胞におけるそれぞれの細胞の各々について、第1の複数の細胞が、20個以上の細胞を含み、複数の注釈付きの細胞状態を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、複数の細胞構成要素が、10個以上の細胞構成要素を含み、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含む、得ることを含む、プロセスによって複数の細胞構成要素モジュールにおける細胞構成要素モジュールを識別することを更に含む。したがって、方法は、複数のベクトルにアクセスするか、又はそれらを形成し、複数のベクトルにおけるそれぞれのベクトルの各々が、(i)複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii)対応する複数のエレメントを含み、対応する複数のエレメントにおけるそれぞれのエレメントの各々が、第1の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を表す対応するカウントを有する。複数のベクトルは、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別するために使用される。複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々は、複数の細胞構成要素のサブセットを含み、複数の細胞構成要素モジュールは、(i)複数の候補細胞構成要素モジュール、及び(ii)複数の細胞構成要素、又はその表現によって次元決定された潜在表現で配置され、複数の細胞構成要素モジュールは、10を超える細胞構成要素モジュールを含む。

30

40

#### 【0338】

1つ以上の第2のデータセットは、電子形式で得られ、1つ以上の第2のデータセットは、第2の複数の細胞におけるそれぞれの細胞の各々について、第2の複数の細胞が、2

50

0 個以上の細胞を含み、目的の生理学的状態を通知する複数の共変量を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含む。したがって、( i ) 第 2 の複数の細胞及び ( i i ) 複数の細胞構成要素又はその表現によって次元決定された細胞構成要素カウントデータ構造が得られる。活性化データ構造は、複数の細胞構成要素又はその表現を共通次元として使用して、細胞構成要素カウントデータ構造及び潜在表現を組み合わせることによって形成され、活性化データ構造は、複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、第 2 の複数の細胞における細胞の各々について、それぞれの活性化重みを含む。

#### 【 0 3 3 9 】

候補細胞構成要素モデルは、( i ) 活性化データ構造を候補モデルに入力したときに、活性化データ構造内に表される細胞構成要素モジュールの各々における複数の共変量における各共変量の不在又は存在の予測と、( i i ) 細胞構成要素モジュールの各々における各共変量の実際の不在又は存在との間の差を使用して訓練され、訓練は、差に応答して、候補細胞構成要素モデルと関連付けられた複数の共変量重みを調整し、複数の共変量重みは、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの共変量の各々について、それぞれの共変量が、活性化データ構造にわたって、それぞれの細胞構成要素モジュールと相関するかどうかを示す対応する重みを含む。候補細胞構成要素モデルを訓練する際に、複数の共変量重みを使用して、複数の候補細胞構成要素モジュールにおける細胞構成要素モジュール（例えば、目的の生理学的状態

#### 【 0 3 4 0 】

いくつかの実施形態において、目的の生理学的状態は、疾患であり、第 1 の複数の細胞が、複数の注釈付きの細胞状態によって立証されるように、疾患を代表する細胞、及び疾患を代表しない細胞を含む。いくつかの実施形態において、複数の注釈付きの細胞状態における注釈付きの細胞状態は、曝露条件下での化合物への第 1 の複数の細胞における細胞の曝露である。いくつかの実施形態において、曝露条件は、曝露期間、化合物の濃度、又は曝露期間及び化合物の濃度の組み合わせである。

#### 【 0 3 4 1 】

いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである。いくつかの実施形態において、第 1 又は第 2 の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動 ( F R E T ) 測定によって決定される。いくつかの実施形態において、第 1 又は第 2 の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、単一細胞リボ核酸 ( R N A ) 配列決定 ( s c R N A - s e q )、s c T a g - s e q、配列決定を使用したトランスポザーゼ - アクセス可能なクロマチンのための単一細胞アッセイ ( s c A T A C - s e q )、C y T O F / S C o P、E - M S / A b s e q、m i R N A - s e q、C I T E - s e q、及びそれらの任意の組み合わせによって決定される。いくつかの実施形態において、複数の細胞構成要素は、1 0 0 ~ 8 , 0 0 0 個の細胞構成要素からなる。

#### 【 0 3 4 2 】

いくつかの実施形態において、複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することは、複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、複数のベクトルに相関モデルを適用することを含む。いくつかの実施形態において、相関モデルは、グラフクラスタリングを含む。いくつかの実施形態において、グラフクラスタリング方法は、ピアソン相関ベースの距離メトリック上のライデン ( L e i d e n ) クラスタリングであるか、又はルーバン ( L o u v a i n ) クラスタリングである。

10

20

30

40

50

## 【0343】

いくつかの実施形態において、複数の細胞構成要素モジュールは、10～2000個の細胞構成要素モジュールからなる。いくつかの実施形態において、複数の構成要素モジュールにおける候補細胞構成要素モジュールの各々は、200～300個の細胞構成要素からなる。

## 【0344】

いくつかの実施形態において、複数の共変量は、細胞バッチ、細胞ドナー、細胞型、疾患状態、又は化学化合物への曝露を含む。

## 【0345】

いくつかの実施形態において、候補細胞構成要素モデルを訓練することは、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する。

10

## 【0346】

ブロック806を参照すると、方法は、複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、(ii)細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。

20

## 【0347】

いくつかの実施形態において、1つ以上の計算された活性化スコアにおける活性化スコアは、それぞれの化合物に対応するそれぞれの細胞構成要素モジュールについてのそれぞれの活性化重みである。例えば、いくつかの実施形態において、活性化スコアは、図2A～図2Bに記載され、図5の活性化データ構造に示されるように得られた活性化重みであり、活性化スコアは、それぞれの(例えば、第1の)細胞構成要素モジュールの活性化(例えば、誘導及び/又は差次的発現)を示し、それぞれの化合物による処置に相関及び/又は応答している。

30

## 【0348】

「モデルアーキテクチャ」と題された上記のセクションに開示されるものなどのモデルの任意の好適な実施形態が企図され、当業者には明らかであろうように、それらの任意の置換、修飾、追加、削除、及び/又は組み合わせが企図される。例えば、いくつかの実施形態において、訓練されたモデルは、ニューラルネットワークを含む。いくつかの実施形態において、ニューラルネットワークは、ReLU活性化を有する完全に接続されたニューラルネットワークである。いくつかの実施形態において、ニューラルネットワークは、メッセージパッシングニューラルネットワークである。いくつかの実施形態において、訓練されたモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

40

## 【0349】

いくつかの実施形態において、訓練されたモデルは、複数のコンポーネントモデルのアンサンブルモデルであり、それぞれの計算された活性化スコアは、複数のコンポーネントモデルにおけるコンポーネントモデルの各々の出力の中心傾向の測定値である。いくつかの実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。いくつかの実施形態

50

において、複数のコンポーネントモデルは、複数のニューラルネットワークを含む。いくつかの実施形態において、複数のニューラルネットワークにおける第1のニューラルネットワークは、ReLU活性化を伴う完全に接続されたニューラルネットワークであり、複数のニューラルネットワークにおける第2のニューラルネットワークは、メッセージパッシングニューラルネットワークである。

**【0350】**

ブロック808を参照すると、訓練は、差に応答して訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。

10

**【0351】**

いくつかの実施形態において、モデルへの入力、複数の活性化スコアを含み、それぞれの活性化スコアの各々は、複数の化合物における化合物の各々について、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールに対応する。それぞれの化合物の各々についてのそれぞれの細胞構成要素モジュールの各々に対応する活性化スコアは、モジュールと化合物との間の関連性（例えば、重み及び/又は相関）を識別するためにマルチタスクモデルを訓練するための標識（例えば、モジュールと化合物との間の関連性の実際の存在又は不在を示す数値的活性化スコア）として機能する。例えば、上述のように、いくつかの実施形態において、複数の細胞構成要素モジュールの第1のサブセットは、目的の生理学的状態に関連付けられ、複数の細胞構成要素モジュールの第2のサブセットは、目的の生理学的状態に関連付けられていない。したがって、いくつかのそのような実施形態において、関連性の実際の存在は、複数の細胞構成要素モジュールの第1のサブセットを標識として使用して訓練データセットに含めることができ、関連性の実際の不在は、複数の細胞構成要素モジュールの第2のサブセットを標識として使用して訓練データセットに含めることができる。

20

**【0352】**

いくつかの実施形態において、複数の化合物は、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも800個、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、少なくとも5000個、少なくとも8000個、少なくとも10,000個、少なくとも20,000個、少なくとも30,000個、少なくとも50,000個、少なくとも80,000個、少なくとも100,000個、少なくとも200,000個、少なくとも500,000個、少なくとも800,000個、少なくとも100万個、又は少なくとも200万個の化合物を含み、複数の化合物における化合物の各々について、モデルへの入力、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの活性化スコアを含む。

30

**【0353】**

いくつかの実施形態において、複数の化合物は、1000万個以下、500万個以下、100万個以下、500,000個以下、100,000個以下、50,000個以下、10,000個以下、8000個以下、5000個以下、2000個以下、1000個以下、800個以下、500個以下、200個以下、又は100個以下の化合物を含み、複数の化合物における化合物の各々について、モデルへの入力、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの活性化スコアを含む。いくつかの実施形態において、複数の化合物は、10~500個、100~10,000個、5000~200,000個、又は10,000~100万個の化合物からなり、複数の化合物における化合物の各々について、モデルへの入力、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの活性化スコアを含む。

40

50

## 【0354】

いくつかの実施形態において、上述したように、複数の数値的活性化スコアにおけるそれぞれの数値的活性化スコアは、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、複数の化合物における化合物の各々についての活性化重みである（例えば、図5の活性化データ構造に示される）。

## 【0355】

上記のように、いくつかの実施形態において、モデルの出力は、複数の化合物におけるそれぞれの化合物（例えば、試験化学化合物）が、複数の細胞構成要素モジュールにおけるそれぞれの1つ以上の細胞構成要素モジュールと関連するかどうかを示す1つ以上の計算された活性化スコアを含む。

## 【0356】

一般に、モデル（例えば、ニューラルネットワーク）を訓練することは、逆伝搬（例えば、勾配降下）を通してそれぞれのモデルについての複数のパラメータ（例えば、重み）を更新することを含む。第一に、入力データ（例えば、複数のモジュールにおけるそれぞれの細胞構成要素モジュールの各々について、複数の化合物におけるそれぞれの化合物の各々についての複数の活性化スコア）がニューラルネットワークに受け入れられ、選択された活性化関数及びパラメータの初期セット（例えば、重み及び/又はハイパーパラメータ）に基づいて出力が計算される、順方向伝搬が実施される。いくつかの実施形態において、パラメータ（例えば、重み及び/又はハイパーパラメータ）は、訓練されていないか、又は部分的に訓練されたモデルに対してランダムに割り当てられる（例えば、初期化される）。いくつかの実施形態において、パラメータは、以前に保存された複数のパラメータから、又は事前に訓練されたモデルから（例えば、転移学習によって）転送される。

## 【0357】

次いで、後方パスが、各層におけるそれぞれのユニットの各々に対応するそれぞれのパラメータの各々についての誤差勾配を計算することによって実施され、各パラメータについての誤差は、ネットワーク出力（例えば、計算された活性化スコアとしてのそれぞれの化合物とそれぞれの細胞構成要素モジュールとの間の関連性の予測された不在又は存在）及び入力データ（例えば、期待値又は真の標識、数値的活性化スコアとしてのそれぞれの化合物とそれぞれの細胞構成要素モジュールとの間の関連性の実際の不在又は存在）に基づいて損失（例えば、誤差）を計算することによって決定される。次いで、パラメータ（例えば、重み）は、計算された損失に基づいて値を調整することによって更新され、それによってモデルを訓練する。

## 【0358】

例えば、機械学習のいくつかの一般的な実施形態において、逆伝搬は、複数の重み（例えば、埋め込み）を含む隠れ層を有するネットワークを訓練する方法である。訓練されていないモデルの出力（例えば、計算された活性化スコアとしての関連性の予測された不在又は存在）は、最初に任意に選択された初期重みのセットを使用して生成される。次いで、（例えば、損失関数を使用して）誤差を計算するために誤差関数を評価することによって、出力を元の入力（例えば、数値的活性化スコアとしての関連性の実際の不在又は存在）と比較する。次いで、重みは、（例えば損失関数に従って）誤差が最小化されるように更新される。いくつかの実施形態において、当業者には明らかであろうように、様々な逆伝搬アルゴリズム及び/又は方法のいずれか1つが、複数の重みを更新するために使用される。

## 【0359】

いくつかの実施形態において、損失関数は、平均平方誤差、二次損失、平均絶対誤差、平均バイアス誤差、ヒンジ、マルチクラスサポートベクトルマシン、及び/又は交差エントロピーである。いくつかの実施形態において、訓練されていないか、又は部分的に訓練されたモデルを訓練することは、勾配降下アルゴリズム及び/又は最小化関数に従って誤差を計算することを含む。いくつかの実施形態において、訓練されていないか、又は部分的に訓練されたモデルを訓練することは、複数の損失関数を使用して複数の誤差を計算す

10

20

30

40

50



ることを含む。いくつかの実施形態において、複数の損失関数における損失関数の各々は、同じ又は異なる重み付け係数を受け取る。

#### 【0360】

図6は、本開示のいくつかの実施形態による、モデルを訓練するための方法の例を示す。活性化データ構造(上部パネル)は、複数のK細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々と、複数のG細胞における細胞の各々との間の関連性を示す複数の活性化スコアを含むモデルへの入力を提供し、細胞の各々は、複数の化合物におけるそれぞれの化合物を表す。複数の細胞構成要素モジュール(中央パネル)におけるそれぞれの細胞構成要素モジュールの各々について、訓練の前に、複数の細胞(例えば、W化合物)によって集合的に表される複数の化合物におけるそれぞれの化合物の各々について、対応する重みを初期化する(例えば、ランダムな重みにする)。したがって、複数の化合物重みは、化合物重みマトリックス(中央パネル)を含む。複数の化合物重みの調整は、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコア(例えば、予測)と、(ii)細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれの数値的活性化スコア(例えば、実際)との間の差を使用して実施される(下部パネル)。いくつかの実施形態において、実際の活性化は、例えば、図2A~図2B及び図14A~図14Dを参照して、以下の「細胞構成要素モジュールの識別」と題されるセクションに記載の細胞構成要素モジュールを識別するための方法を使用して得られ、複数の共変量は複数の化合物を含む。次いで、訓練(例えば、化合物重みの調整)は、訓練されたモデルが形成されるまで(例えば、最小数の調整の完了及び/又は最小性能閾値の達成によって)実施することができる。

10

20

#### 【0361】

いくつかの実施形態において、誤差関数は、計算された損失に比例する量によって1つ以上のパラメータの値を調整することによって、モデル(例えば、ニューラルネットワーク)における1つ以上のパラメータ(例えば、重み)を更新するために使用され、それによってモデルを訓練する。いくつかの実施形態において、パラメータが調整される量は、パラメータが更新される程度又は重大度(例えば、より小さい又はより大きい調整)を指示する学習率ハイパーパラメータによって計測される。したがって、いくつかの実施形態において、訓練は、学習率に基づいて、複数のパラメータの全て又はサブセットを更新する。いくつかの実施形態において、学習率は、差次的学習率である。

30

#### 【0362】

いくつかの実施形態において、モデル(例えば、ニューラルネットワーク)を訓練することは、対応する複数の隠れニューロンにおける隠れニューロンの各々の対応するパラメータに対する正規化を更に使用する。例えば、いくつかの実施形態において、正規化は、損失関数にペナルティを追加することによって実施され、ペナルティは、ニューラルネットワークにおけるパラメータの値に比例する。一般に、正規化は、1つ以上のパラメータにペナルティを追加することによってモデルの複雑性を低減し、それらのパラメータと関連付けられたそれぞれの隠れニューロンの重要性を低下させる。そのような実践は、より一般化されたモデルをもたらす、データの過剰適合を低減することができる。いくつかの実施形態において、正規化は、L1又はL2ペナルティを含む。例えば、いくつかの好ましい実施形態において、正規化は、より低い及びより高いパラメータに対するL2ペナルティを含む。いくつかの実施形態において、正規化は、空間正規化(例えば、先験的及び/又は実験的知識に基づいて決定される)又はドロップアウト正規化を含む。いくつかの実施形態において、正規化は、独立して最適化されるペナルティを含む。

40

#### 【0363】

いくつかの実施形態において、モデルに関連付けられた複数の化合物重みを調整すること(例えば、予測された標識と実際の標識との間の差に回答する)を含む訓練プロセスは、複数の訓練インスタンスにおける訓練インスタンスの各々に対して繰り返される。

50

## 【0364】

いくつかの実施形態において、複数の訓練インスタンスは、少なくとも3、少なくとも4、少なくとも5、少なくとも6、少なくとも7、少なくとも8、少なくとも9、少なくとも10、少なくとも50、少なくとも100、少なくとも500、少なくとも1000、少なくとも2000、少なくとも3000、少なくとも4000、少なくとも5000、又は少なくとも7500の訓練インスタンスを含む。いくつかの実施形態において、複数の訓練インスタンスは、10, 000以下、5000以下、1000以下、500以下、100以下、又は50以下の訓練インスタンスを含む。いくつかの実施形態において、複数の訓練インスタンスは、3~10、5~100、100~5000、又は1000~10,000の訓練インスタンスを含む。いくつかの実施形態において、複数の訓練インスタンスは、3以上の訓練インスタンスから始まり、10,000以下の訓練インスタンスで終わる別の範囲内にある。

10

## 【0365】

いくつかのそのような実施形態において、訓練は、複数の訓練インスタンスにわたって（例えば、逆伝搬を介して）モデルのパラメータの調整を繰り返すことを含み、したがって、それぞれの化合物がそれぞれの細胞構成要素モジュールと相関するかどうかを示す際のモデルの精度を増加させる。

## 【0366】

いくつかの実施形態において、訓練は、転移学習を含む。転移学習は、例えば、定義のセクションに更に記載されている（上記の「訓練されていないモデル」を参照されたい）。

20

## 【0367】

いくつかの実施形態において、訓練されていないか、又は部分的に訓練されたモデルを訓練することは、誤差関数の第1の評価に続いて訓練されたモデルを形成する。いくつかのそのような実施形態において、訓練されたモデルは、誤差関数の第1の評価に基づいて、1つ以上のパラメータの第1の更新に続いて形成される。いくつかの代替の実施形態において、訓練されたモデルは、誤差関数の少なくとも1回、少なくとも2回、少なくとも3回、少なくとも4回、少なくとも5回、少なくとも6回、少なくとも7回、少なくとも8回、少なくとも9回、少なくとも10回、少なくとも20回、少なくとも30回、少なくとも40回、少なくとも50回、少なくとも100回、少なくとも500回、少なくとも1000回、少なくとも10,000回、少なくとも50,000回、少なくとも100,000回、少なくとも200,000回、少なくとも500,000回、又は少なくとも100万回の評価に続いて形成される。いくつかのそのような実施形態において、訓練されたモデルは、誤差関数の少なくとも1回、少なくとも2回、少なくとも3回、少なくとも4回、少なくとも5回、少なくとも6回、少なくとも7回、少なくとも8回、少なくとも9回、少なくとも10回、少なくとも20回、少なくとも30回、少なくとも40回、少なくとも50回、少なくとも100回、少なくとも500回、少なくとも1000回、少なくとも10,000回、少なくとも50,000回、少なくとも100,000回、少なくとも200,000回、少なくとも500,000回、又は少なくとも100万回の評価に基づいて、1つ以上のパラメータの少なくとも1回、少なくとも2回、少なくとも3回、少なくとも4回、少なくとも5回、少なくとも6回、少なくとも7回、少なくとも8回、少なくとも9回、少なくとも10回、少なくとも20回、少なくとも30回、少なくとも40回、少なくとも50回、少なくとも100回、少なくとも500回、少なくとも1000回、少なくとも10,000回、少なくとも50,000回、少なくとも100,000回、少なくとも200,000回、少なくとも500,000回、又は少なくとも100万回の更新に続いて形成される。

30

40

## 【0368】

いくつかの実施形態において、訓練されたモデルは、モデルが最小性能要件を満たす場合に形成される。例えば、いくつかの実施形態において、訓練されたモデルについて計算された誤差が、誤差関数の評価（例えば、各化合物と各細胞構成要素モジュールとの間の

50

予測された関連性と実際の関連性との間の差)に続いて、誤差閾値を満たす場合に、訓練されたモデルが形成される。いくつかの実施形態において、誤差関数によって計算される誤差は、誤差が20パーセント未満、18パーセント未満、15パーセント未満、10パーセント未満、5パーセント未満、又は3パーセント未満である場合に誤差閾値を満たす。

#### 【0369】

例示的な実施形態において、モデルを訓練することは、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する。

10

#### 【0370】

いくつかの実施形態において、訓練することは、回帰モデルに従って、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々についてのそれぞれの化合物の各々と関連付けられた差の各々に応答して、訓練されていないモデルと関連付けられた複数のパラメータを調整する。いくつかの実施形態において、回帰モデルは、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々についてのそれぞれの化合物の各々と関連付けられた差の各々の最小二乗誤差を最適化する。

#### 【0371】

モデル訓練の前述の説明は、化合物と細胞構成要素モジュールとの間の関連性を示す活性化スコアを得て、使用することを記載しているが、実際には、化合物と目的の任意の他の生理学的状態、又はその任意の細胞プロセスとの間の関連性を示す活性化スコアは、化合物を生理学的状態と関連付けるための訓練及びモデルの使用において企図されている。例えば、以下のセクションで説明されるように、本開示の別の態様は、摂動シグネチャを使用してモデルを訓練することを含む。具体的には、いくつかの実施形態において、モデルは、訓練標識として摂動シグネチャのための数値的活性化スコアを使用して訓練される。次いで、訓練されたモデルは、方法700において記載されるように、出力として、モデルへの化学構造フィンガープリントの入力に応答して、1つ以上の計算された活性化スコアを得るために使用され、1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、摂動シグネチャのセットにおける対応する摂動シグネチャを表す。

20

30

#### 【0372】

摂動シグネチャについての数値的活性化スコアの取得。

したがって、本開示の別の態様は、化学化合物を目的の生理学的状態と関連付けるための方法900を提供する。いくつかの実施形態において、目的の生理学的状態は、疾患である。

#### 【0373】

ブロック902を参照すると、方法は、複数の化合物における化合物の各々の化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得る。「生理学的状態」及び「化合物」と題する上記のセクションに開示されるような、生理学的状態、化合物、フィンガープリント、及び/又はフィンガープリントを得る方法の任意の好適な実施形態は、当業者には明白であろうように、それらの任意の置換、修飾、追加、欠失、及び/又は組み合わせを含むことが企図される。

40

#### 【0374】

例えば、いくつかの実施形態において、複数の化合物は、 $10 \sim 1 \times 10^6$ 個の化合物である。いくつかの実施形態において、複数の化合物は、 $100 \sim 100,000$ 個の化合物である。いくつかの実施形態において、複数の化合物は、 $1000 \sim 100,000$ 個の化合物である。いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、2000ダルトン未満の分子量を有する有機化合物である。いくつかの実施形態において、複数の化学化合物における化学化合物の各々は、5つの基準のリピンスキーの法則の各々を満たす。いくつかの実施形態において、複数の化学化合物における化学化

50

化合物の各々は、5つの基準のリピンスキーの法則のうち少なくとも3つの基準を満たす。いくつかの実施形態において、それぞれのフィンガープリントの各々は、SMILES Transformer、ECFP4、RNNS2S、又はGraphConvを使用して、化学構造から生成される。

#### 【0375】

ブロック904を参照すると、方法は、複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子形式で得ることを含む。「摂動シグネチャ」と題された上記のセクションに開示される摂動シグネチャの任意の好適な実施形態は、当業者に明白であるように、それらの任意の置換、修正、追加、欠失、及び/又は組み合わせを含むことが企図される。

10

#### 【0376】

例えば、いくつかの実施形態において、摂動シグネチャのセットは、単一の摂動シグネチャである。いくつかの実施形態において、摂動シグネチャのセットは、複数の摂動シグネチャである。いくつかの実施形態において、摂動シグネチャのセットは、200~500個の摂動シグネチャである。いくつかの実施形態において、摂動シグネチャのセットは、5つ以上の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、10個以上の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、100個以上の摂動シグネチャを含む。いくつかの実施形態において、摂動シグネチャのセットは、複数の摂動シグネチャであり、複数の摂動シグネチャの第1のサブセットは、目的の生理学的状態と関連付けられ、複数の摂動シグネチャの第2のサブセットは、目的の生理学的状態と関連付けられていない。

20

#### 【0377】

ブロック906を参照すると、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々が、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうち一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。

30

#### 【0378】

いくつかの実施形態において、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアは、変化していない細胞状態と、変化した細胞状態との間の差次的な細胞構成要素存在量の測定値を表す単一細胞遷移シグネチャに電子形式でアクセスすることを含む手順によって得られる。変化した細胞状態は、変化していない細胞状態から変化した細胞状態への細胞遷移を通じて生じ、(i)変化していない細胞状態、(ii)変化した細胞状態、及び(iii)変化していない細胞状態から変化した細胞状態への遷移のうち少なくとも1つは、目的の生理学的状態と関連付けられる。単一細胞遷移シグネチャは、参照の複数の細胞構成要素の識別と、複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、変化していない細胞状態と変化した細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する第1の有意性スコアとを含む。それぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定するために、単一細胞遷移シグネチャ及びそれぞれの摂動シグネチャを比較する。

40

#### 【0379】

いくつかの実施形態において、単一細胞遷移シグネチャと摂動シグネチャとを比較して、それぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定することは、単一細胞遷移シグネチャの参照の複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの摂動シグネチャにおける対応する細胞構成要素の対応する有意性スコ

50

アに対する単一細胞遷移シグネチャにおけるそれぞれの細胞構成要素の第 1 の有意性スコアを比較することを含む。

【 0 3 8 0 】

いくつかの実施形態において、単一細胞遷移シグネチャと摂動シグネチャとを比較して、それぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定することは、単一細胞遷移シグネチャの参照の複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの摂動シグネチャにおける複数の細胞構成要素における対応する細胞構成要素の各々の対応する有意性スコアに対する単一細胞遷移シグネチャにおける複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々の有意性スコアを比較することを含む。

10

【 0 3 8 1 】

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、摂動シグネチャのセットにおける他の摂動シグネチャと比較して、単一細胞遷移シグネチャに対するそれぞれの摂動シグネチャの関連性の相対的なランキングである。いくつかの実施形態において、相対的なランキングは、ウィルコクソンの順位和検定、t 検定、ロジスティック回帰、又は一般化線形モデルによって決定される。いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、ランキングに基づいていない。

【 0 3 8 2 】

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、それぞれの摂動シグネチャについてのそれぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての、対応する有意性スコアの中心傾向の測定値である。いくつかの実施形態において、中心傾向の測定値は、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての対応する有意性スコアの算術平均、加重平均、ミッドレンジ、ミッドヒンジ、トリミアン、ウィンザライズド平均、平均、又はモードである。

20

【 0 3 8 3 】

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、( i ) それぞれの摂動シグネチャについての、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての、対応する有意性スコアの中心傾向の測定値と、( i i ) 単一細胞遷移シグネチャについての、複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々についての、対応する第 1 の有意性スコアの中心傾向の測定値との間の差である。

30

【 0 3 8 4 】

いくつかの実施形態において、単一細胞遷移シグネチャの変化していない細胞状態が、それぞれの摂動シグネチャの第 1 の細胞状態又は第 2 の細胞状態と同じである。いくつかの実施形態において、単一細胞遷移シグネチャの変化していない細胞状態が、それぞれの摂動シグネチャの第 1 の細胞状態及び第 2 の細胞状態の両方とは異なる。

【 0 3 8 5 】

いくつかの実施形態において、方法は、単一細胞遷移シグネチャの参照の複数の細胞構成要素、及びそれぞれの摂動シグネチャのそれぞれの複数の細胞構成要素を剪定して、転写因子と比較することを制限することを更に含む。いくつかの実施形態において、方法は、単一細胞遷移シグネチャの参照の複数の細胞構成要素、及びそれぞれの摂動シグネチャのそれぞれの複数の細胞構成要素を剪定して、別の細胞構成要素の種類（例えば、遺伝子、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、及び / 又はそれらの組み合わせ）との比較を制限することを更に含む。いくつかの実施形態において、参照の複数の細胞構成要素及びそれぞれの複数の細胞構成要素は、剪定されない。

40

【 0 3 8 6 】

いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動状態は、複数の化合物における化合物に曝露されていない対照細胞によって表される。いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動状態は、それぞれの摂動シグネチャに関連付けられた化合物以外の複数の化

50

学化合物における化学化合物に曝露されている無関係の摂動細胞にわたる平均によって表される。

【0387】

上述のように、いくつかの実施形態において、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャは、非限定的な例として、参照により本明細書に組み込まれる、2019年7月15日に出願された「Methods of Analyzing Cells」と題された米国特許出願第16/511,691号に開示された方法のいずれかを使用して決定され得る。

【0388】

それぞれの摂動シグネチャは、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含む。それぞれの第1の細胞状態及び第2の細胞状態のうち的一方は、非摂動細胞状態であり、他方は、それぞれの摂動シグネチャに対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。更に、上述したように、それぞれの摂動シグネチャは、数値的活性化スコアを含む。いくつかの実施形態において、それぞれの摂動シグネチャについての数値的活性化スコアは、連続スケール上の絶対値である。いくつかの実施形態において、それぞれの摂動シグネチャについての数値的活性化スコアは、以下でより詳細に説明されるように、相対的なランキングである。

【0389】

いくつかの実施形態において、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアは、変化していない細胞状態と、変化した細胞状態との間の差次的な細胞構成要素存在量の測定値を表す単一細胞遷移シグネチャに電子形式でアクセスすることを含む手順によって得られる。ここで、変化した細胞状態は、変化していない細胞状態から変化した細胞状態への細胞遷移を通して発生する。更に、(i)変化していない細胞状態、(ii)変化した細胞状態、及び(iii)変化していない細胞状態から変化した細胞状態への遷移のうち少なくとも1つが、目的の生理学的状態と関連付けられる。

【0390】

単一細胞遷移シグネチャは、参照の複数の細胞構成要素の識別と、複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、変化していない細胞状態と変化した細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する第1の有意性スコアとを含む。いくつかの実施形態において、単一細胞遷移シグネチャは、参照により本明細書に組み込まれる、2019年7月15日に出願された「Methods of Analyzing Cells」と題される米国特許出願第16/511,691号に開示される方法のいずれかを使用して決定される。

【0391】

一度得られると、単一細胞遷移シグネチャは、それぞれの摂動シグネチャと比較され、それによってそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定する。いくつかの実施形態において、2019年7月15日に出願された「Methods of Analyzing Cells」と題された米国特許出願第16/511,691号に開示された、単一細胞遷移シグネチャをそれぞれの摂動シグネチャと比較して、複数の摂動シグネチャにおける他の摂動シグネチャに対して、それぞれの摂動シグネチャの相対的なランキングを決定するための方法のいずれかを使用することができ、例えば、そのような相対的なランキングは、次いで、それぞれの摂動シグネチャのそれぞれの数値的活性化スコアとみなされるであろう。

【0392】

いくつかの実施形態において、単一細胞遷移シグネチャ及び摂動シグネチャを比較して

、それぞれの摂動シグネチャのそれぞれの数値的活性化スコアを決定することは、単一細胞遷移シグネチャの参照の複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の第1の有意性スコアを、それぞれの摂動シグネチャにおける対応する細胞構成要素の対応する有意性スコアと比較することを含む。いくつかのそのような実施形態において、それぞれの摂動シグネチャの活性化スコアは、摂動シグネチャのセットにおける他の摂動シグネチャと比較して、単一細胞遷移シグネチャに対するそれぞれの摂動シグネチャの関連性の相対的なランキングである。いくつかのそのような実施形態において、相対的なランキングは、ウィルコクソンの順位和検定、t検定、ロジスティック回帰、又は一般化線形モデルによって決定される。いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、それぞれの摂動シグネチャの関連性の相対的なランキングではなく、むしろ、単一細胞遷移シグネチャに対する他の摂動シグネチャのランキングとは独立して決定される。

10

**【0393】**

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、ランキングに基づいていない。例えば、いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、それぞれの摂動シグネチャについてのそれぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての対応する有意性スコアを含む複数の有意性スコアである。

**【0394】**

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、それぞれの摂動シグネチャについてのそれぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての、対応する有意性スコアの中心傾向の測定値である。いくつかの実施形態において、中心傾向の測定値は、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての対応する有意性スコアの算術平均、加重平均、ミッドレンジ、ミッドヒンジ、トリミアン、ウィンザライズド平均、平均、又はモードである。

20

**【0395】**

いくつかの実施形態において、それぞれの摂動シグネチャの活性化スコアは、(i)それぞれの摂動シグネチャについての、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々についての、対応する有意性スコアの中心傾向の測定値と、(ii)単一細胞遷移シグネチャについての、複数の参照細胞構成要素におけるそれぞれの細胞構成要素の各々についての、対応する第1の有意性スコアの中心傾向の測定値との間の差である。

30

**【0396】**

一実施形態において、単一細胞遷移シグネチャとそれぞれの摂動シグネチャとの間の比較を実施するために、摂動シグネチャの細胞構成要素は、マトリックスとして表される。マトリックスの各行は、単一の摂動(例えば、複数の化合物における単一の化合物)と関連付けられる。マトリックス上の各列は、それぞれの状態間の差次的な存在量を示す細胞構成要素のうちの一つと関連付けられる。マトリックスの各エントリーは、特定の摂動シグネチャについて識別された細胞構成要素についての有意性スコア(例えば、p値、tスコア)を含む。このマトリックスは、単一細胞遷移シグネチャにある細胞構成のみを含むようにフィルタリングされる。このフィルタリングは、閾値p値、細胞構成要素の閾値数の使用などを使用して達成され得る。

40

**【0397】**

マトリックス内の各有意性スコアは、個別のマッチングスコアと置き換えられる。各有意性スコアを個別のマッチングスコアと置き換えるために、細胞遷移についての有意に上方制御された細胞構成要素及び細胞遷移についての有意に下方制御された細胞構成要素を識別する。単一細胞遷移シグネチャによって識別される有意に上方制御された細胞構成要素の各々について、細胞構成要素がその摂動(例えば、化学組成物)についての摂動シグネチャについても有意に上方制御されている場合、その細胞構成要素/摂動の組み合わせについてのマトリックスにおける有意性スコアは、「1」の個別のマッチングスコアと置

50

き換えられる。細胞構成要素が単一細胞遷移シグネチャと比較して摂動シグネチャに対して有意に下方制御されている場合、その細胞構成要素 / 摂動の組み合わせについてのマトリックスにおける有意性スコアは、「 - 2 」の個別のマッチングスコアと置き換えられる。細胞構成要素が摂動シグネチャに対して有意に上方制御又は下方制御されていない場合、細胞構成要素 / 摂動の組み合わせについてのマトリックスにおける有意性スコアは、「 0 」の個別のマッチングスコアと置き換えられる。

**【 0 3 9 8 】**

逆に、単一細胞遷移シグネチャにおいて識別された有意に下方制御された細胞構成要素の各々について、細胞構成要素が摂動についても有意に下方制御されている場合、その細胞構成要素 / 摂動の組み合わせについてのマトリックスにおける有意性スコアは、「 - 1 」の個別のマッチングスコアと置き換えられる。細胞構成要素が摂動に対して有意に上方制御されている場合、その細胞構成要素 / 摂動の組み合わせについてのマトリックスにおける有意性スコアは、「 2 」の個別のマッチングスコアと置き換えられる。細胞構成要素が摂動シグネチャに対して有意に上方制御又は下方制御されていない場合、その細胞構成要素 / 摂動の組み合わせについてのマトリックスにおける有意性スコアは、「 0 」の個別のマッチングスコアと置き換えられる。当業者は、いくつかの実施形態において、これらの特定のスコア置換が他の数値で置換され得ることを理解するであろう。更に、上方制御又は下方制御の代わりに、細胞構成要素の各々についての閾値存在量値の使用が使用され得、次いで、所与の細胞構成要素が閾値存在量値を上回るか、又は下回るかどうかの考慮が、前述のクラス標識（例えば、「 - 1 」、「 2 」、「 0 」など）をマトリックスの各エレメントに割り当てる際に行われる。

**【 0 3 9 9 】**

結果は、摂動の数（複数の化学組成物における化学組成物の数、したがって複数の摂動シグネチャにおける摂動シグネチャの数）によって与えられる行の数と、上記のマトリックスエレメントエントリーがマッチングスコアを表す単一細胞遷移からの差次的細胞構成要素によって与えられる列の数とのマトリックスである。

**【 0 4 0 0 】**

上記のように、マトリックス内の有意性スコアを個別のマッチングスコアに置き換えた後、マトリックスの各行における個別のマッチングスコアを合計して、各行についての合計されたマッチングスコアを生成する。次いで、各々が摂動シグネチャに対応するマトリックスの行は、合計したマッチングスコアを減少させる順序でランク付けされる。最上位の行は、単一細胞遷移シグネチャの識別された細胞遷移と関連付けられる可能性が最も高い摂動シグネチャと関連付けられる。更に、行の各々のランキングは、行の各々に対応する摂動シグネチャについての活性化スコアとして使用することができる。

**【 0 4 0 1 】**

いくつかの実施形態において、マトリックスにおける各行の合計マッチングスコアについて、偽の細胞構成要素発見率の推定は、参照により本明細書に組み込まれる、「 Methods of Analyzing Cells 」と題された米国特許出願第 1 6 / 5 1 1 , 6 9 1 号で説明されるように推定される。

**【 0 4 0 2 】**

ある特定の実施形態において、摂動（例えば、特定の化学組成物への細胞の曝露）の共変量が存在し得る。例えば、化学組成物の共変量は、化学組成物の特定の用量、化学組成物に曝露された細胞が細胞構成要素を定量化するために測定される時間、及び / 又は化学組成物に曝露された細胞の同一性（例えば、細胞株）を含み得る。いくつかの実施形態において、摂動（例えば、特定の化学組成物への細胞の曝露）は、その共変量の閾値量も特定の細胞遷移に影響すると予測される場合にのみ、特定の細胞遷移に影響すると予測される。言い換えれば、いくつかの実施形態において、特定の摂動シグネチャの数値的活性化スコアは、特定の摂動シグネチャの化学組成物の共変量が、単一細胞遷移スコアと関連付けられた特定の細胞遷移にも影響を与えると予測されるかどうかによって少なくとも部分的に決定される。

10

20

30

40

50



## 【0403】

それぞれの摂動シグネチャを単一細胞遷移シグネチャと比較する代替方法を使用して、それぞれの摂動シグネチャの数値的活性化スコアを決定し得る。例えば、細胞構成要素は、ウェブインターフェースを使用してデータベースに適合され得る（例えば、amp.ppharm.mssm.edu/L1000CDS2/#/indexのワールドワイドウェブ上のL1000CDS2.Anultra-fastLINCSL1000CharacteristicDirectionSignatureSearchEngineなど）。

## 【0404】

いくつかの実施形態において、単一細胞遷移シグネチャの変化していない細胞状態が、それぞれの摂動シグネチャの第1の細胞状態又は第2の細胞状態と同じである。いくつかの実施形態において、単一細胞遷移シグネチャの変化していない細胞状態が、それぞれの摂動シグネチャの第1の細胞状態及び第2の細胞状態の両方とは異なる。

10

## 【0405】

いくつかの実施形態において、方法は、単一細胞遷移シグネチャの参照の複数の細胞構成要素、及びそれぞれの摂動シグネチャのそれぞれの複数の細胞構成要素を剪定して、転写因子と比較することを制限することを更に含む。いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動状態は、複数の化合物における化合物に曝露されていない対照細胞によって表される。

## 【0406】

いくつかの実施形態において、複数の摂動シグネチャにおけるそれぞれの摂動シグネチャの摂動状態は、それぞれの摂動シグネチャに関連付けられた化合物以外の複数の化学化合物における化学化合物に曝露されている無関係の摂動細胞にわたる平均によって表される。

20

## 【0407】

ブロック908を参照すると、方法は、複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、(ii)摂動シグネチャのセットにおける対応する化合物についてのそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。

30

## 【0408】

「モデルアーキテクチャ」と題された上記のセクションに開示されるものなどのモデルの任意の好適な実施形態が企図され、当業者には明らかであろうように、それらの任意の置換、修飾、追加、削除、及び/又は組み合わせが企図される。例えば、いくつかの実施形態において、訓練されたモデルは、ニューラルネットワークを含む。いくつかの実施形態において、ニューラルネットワークは、ReLU活性化を有する完全に接続されたニューラルネットワークである。いくつかの実施形態において、ニューラルネットワークは、メッセージパッシングニューラルネットワークである。いくつかの実施形態において、訓練されたモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。

40

## 【0409】

いくつかの実施形態において、訓練されたモデルは、複数のコンポーネントモデルのアンサンブルモデルであり、それぞれの計算された活性化スコアは、複数のコンポーネントモデルにおけるコンポーネントモデルの各々の出力の中心傾向の測定値である。いくつかの実施形態において、複数のコンポーネントモデルは、ロジスティック回帰モデル、ニューラルネットワークモデル、サポートベクトルマシンモデル、ナイーブベイズモデル、最

50

近傍モデル、ブーストツリーモデル、ランダムフォレストモデル、決定木モデル、多項ロジスティック回帰モデル、線形モデル、又は線形回帰モデルを含む。いくつかの実施形態において、複数のコンポーネントモデルは、複数のニューラルネットワークを含む。いくつかの実施形態において、複数のニューラルネットワークにおける第1のニューラルネットワークは、ReLU活性化を伴う完全に接続されたニューラルネットワークであり、複数のニューラルネットワークにおける第2のニューラルネットワークは、メッセージパッシングニューラルネットワークである。

【0410】

ブロック910を参照すると、訓練は、差に応答して訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。

10

【0411】

「モデル訓練」と題された上記のセクションに開示されるものなど、訓練されていないか、又は部分的に訓練されたモデルを訓練するための任意の好適な方法及び実施形態は、当業者に明らかであろうように、それらの任意の置換、修飾、追加、欠失、及び/又は組み合わせを含むことが企図される。

【0412】

いくつかの実施形態に関して、モデルへの入力は、複数の活性化スコアを含み、それぞれの活性化スコアの各々は、複数の化合物における化合物の各々について、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャに対応する。それぞれの化合物の各々についてのそれぞれの摂動シグネチャの各々に対応する活性化スコアは、摂動シグネチャと化合物との間の関連性（例えば、重み及び/又は相関）を識別するためにマルチタスクモデルを訓練するための標識（例えば、摂動シグネチャと化合物との間の関連性の実際の存在又は不在を示す数値的活性化スコア）として機能する。例えば、上述のように、いくつかの実施形態において、複数の摂動シグネチャの第1のサブセットは、目的の生理学的状態と関連付けられ、複数の摂動シグネチャの第2のサブセットは、目的の生理学的状態と関連付けられていない。したがって、いくつかのそのような実施形態において、関連性の実際の存在は、複数の摂動シグネチャの第1のサブセットを標識として使用して訓練データセットに含めることができ、関連性の実際の不在は、複数の摂動シグネチャの第2のサブセットを標識として使用して訓練データセットに含めることができる。

20

30

【0413】

いくつかの実施形態において、訓練することは、回帰モデルに従って、摂動シグネチャsのセットにおけるそれぞれの摂動シグネチャの各々についての対応する化合物の各々と関連付けられた差の各々に応答して、訓練されていないモデルと関連付けられた複数のパラメータを調整する。いくつかの実施形態において、回帰モデルは、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々についての対応する化合物の各々と関連付けられた差の各々の最小二乗誤差を最適化する。

【0414】

いくつかの実施形態において、モデルは、細胞構成要素モジュール、摂動シグネチャ、又は両方についての活性化スコアに基づいて、化合物を目的の生理学的状態と関連付けるために訓練及び/又は使用される。いくつかの実施形態において、モデルは、複数のドメイン（例えば、モジュール及び/又は摂動シグネチャなどの標識タイプ）及び/又はデータタイプ（例えば、遺伝子発現プロファイル、メタボロミクス、プロテオミクス、エピジェネティクスなどの分析物及び/又は細胞構成要素）についての活性化スコアに基づいて、化合物を目的の生理学的状態と関連付けるために訓練及び/又は使用される。いくつかの実施形態において、モデルは、任意の1つ以上の目的の生理学的状態（例えば、化合物の毒性、疾患状態の解消など）についての活性化スコアに基づいて、化合物を目的の生理学的状態と関連付けるために訓練及び/又は使用される。いくつかの実施形態において、モデルは、複数のシステムにわたって訓練され、システムは、本明細書に開示される任意

40

50

の1つ以上の生理学的状態、任意の1つ以上のドメイン、及び/若しくは任意の1つ以上のデータタイプ、又は当業者に明白であろう任意の置換、修飾、追加、欠失、及び/又は組み合わせを指す。例えば、いくつかの実施形態において、モデルは、試験化学化合物、毒性の遺伝子モジュール特性の活性化、及び疾患解消を示す摂動シグネチャの間の関連性を集合的に決定するように共同訓練される。

【0415】

追加の実施形態。

本開示の別の態様は、1つ以上のプロセッサ及びメモリを含むコンピュータシステムを提供し、メモリは、試験化学化合物を目的の生理学的状態と関連付けるための方法を実施するための命令を格納する。方法は、試験化学化合物の化学構造のフィンガープリントを得ることと、フィンガープリントをモデルに入力することと、を含み、モデルは、100以上のパラメータを含み、モデルは、フィンガープリントのモデルへの入力にตอบสนองして、1つ以上の計算された活性化スコアを出力し、1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、細胞構成要素モジュールのセットにおける対応する細胞構成要素モジュールを表し、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々は、複数の細胞構成要素の独立したサブセットを含み、細胞構成要素モジュールのセットにおける第1の細胞構成要素モジュールは、目的の生理学的状態と関連付けられる。方法は、第1の細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態と識別することを更に含む。

10

20

【0416】

本開示の別の態様は、試験化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体を提供し、コンピュータは、1つ以上のプロセッサ及びメモリを含み、1つ以上のコンピュータプログラムは、方法を実施するコンピュータによって実行可能な命令を集合的に符号化する。方法は、試験化学化合物の化学構造のフィンガープリントを得ることと、フィンガープリントをモデルに入力することと、を含み、モデルは、100以上のパラメータを含み、モデルは、フィンガープリントのモデルへの入力にตอบสนองして、1つ以上の計算された活性化スコアを出力し、1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、細胞構成要素モジュールのセットにおける対応する細胞構成要素モジュールを表し、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々は、複数の細胞構成要素の独立したサブセットを含み、細胞構成要素モジュールのセットにおける第1の細胞構成要素モジュールは、目的の生理学的状態と関連付けられる。方法は、第1の細胞構成要素モジュールについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態と識別することを更に含む。

30

【0417】

本開示の更に別の態様は、1つ以上のプロセッサ及びメモリを含むコンピュータシステムを提供し、メモリは、試験化学化合物を目的の生理学的状態と関連付けるための方法を実施するための命令を格納する。方法は、試験化学化合物の化学構造のフィンガープリントを得、フィンガープリントをモデルに入力することと、を含み、モデルは100以上のパラメータを含む。モデルは、フィンガープリントのモデルへの入力にตอบสนองして1つ以上の計算された活性化スコアを出力する。1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、摂動シグネチャのセットにおける対応する摂動シグネチャを表す。摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々は、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状

40

50

態のうちの他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。方法は、摂動シグネチャのセットにおける第1の摂動シグネチャについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態と識別することを更に含む。

#### 【0418】

本開示の別の態様は、試験化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体を提供し、コンピュータは、1つ以上のプロセッサ及びメモリを含み、1つ以上のコンピュータプログラムは、方法を実施するコンピュータによって実行可能な命令を集散的に符号化する。方法は、試験化学化合物の化学構造のフィンガープリントを得、フィンガープリントをモデルに入力することを含み、モデルは100以上のパラメータを含む。モデルは、フィンガープリントのモデルへの入力にตอบสนองして1つ以上の計算された活性化スコアを出力する。1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、摂動シグネチャのセットにおける対応する摂動シグネチャを表す。摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々は、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうちの一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうちの他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。方法は、摂動シグネチャのセットにおける第1の摂動シグネチャについてのそれぞれの計算された活性化スコアが、第1の閾値基準を満たす場合、化学化合物を目的の生理学的状態と識別することを更に含む。

#### 【0419】

本開示の更に別の態様は、1つ以上のプロセッサ及びメモリを含むコンピュータシステムを提供し、メモリは化学化合物を目的の生理学的状態と関連付けるための方法を実施するための命令を格納する。方法は、複数の化合物における化合物の各々の化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることを含む。方法は、複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む、得ること、を含む。この方法は、複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、(ii)細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。訓練は、差にตอบสนองして訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。

#### 【0420】

本開示の別の態様は化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体を提供し、コンピュータは、1つ以上のプロセッサ及びメモリを含み、1つ以上のコンピュータプログラムは、方法を実施するコンピュータによって実行可能な命令を集散的に符号化する。方法は、複数の化合物における化合物の各々の化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得る

10

20

30

40

50

ことを含む。方法は、複数の化合物における化合物の各々についての細胞構成要素モジュールのセットにおける細胞構成要素モジュールの各々のそれぞれの数値的活性化スコアを電子形式で得ることであって、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々が、複数の細胞構成要素の独立したサブセットを含む、得ること、を更に含む。この方法は、複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、(ii)細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。訓練は、差に回答して訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。

#### 【0421】

本開示の更に別の態様は、1つ以上のプロセッサ及びメモリを含むコンピュータシステムを提供し、メモリは、化学化合物を目的の生理学的状態と関連付けるための命令を格納し、方法は、複数の化合物における化合物の各々の化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ることを含む。方法は、複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子形式で得ることを更に含む。摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々は、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうち的一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。方法は、複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、(ii)摂動シグネチャのセットにおける対応する化合物についてのそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。訓練は、差に回答して訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。

#### 【0422】

本開示の別の態様は、化学化合物を目的の生理学的状態と関連付けるための、コンピュータによって実行可能な1つ以上のコンピュータプログラムを格納する非一時的なコンピュータ可読媒体を提供し、コンピュータは、1つ以上のプロセッサ及びメモリを含み、1つ以上のコンピュータプログラムは、複数の化合物における化合物の各々の化学構造のそれぞれのフィンガープリントを電子形式で得、それによって複数のフィンガープリントを得ること、を含む、方法を実施するコンピュータによって実行可能な命令を集合的に符号化する。方法は、複数の化合物における対応する化合物の各々についての摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々のそれぞれの数値的活性化スコアを電子形式で得ることを更に含む。摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々は、それぞれの複数の細胞構成要素の識別と、それぞれの複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞構成要素の存在量の変化

と、それぞれの第1の細胞状態とそれぞれの第2の細胞状態との間の細胞状態の変化との間の関連性を定量化する対応する有意性スコアと、を含み、それぞれの第1の細胞状態及び第2の細胞状態のうち一方が、非摂動細胞状態であり、それぞれの第1の細胞状態及び第2の細胞状態のうち他方が、対応する化合物への細胞の曝露によって引き起こされるそれぞれの摂動細胞状態である。方法は、複数の化合物におけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、摂動シグネチャのセットにおけるそれぞれの摂動シグネチャの各々について、(i)それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの摂動シグネチャについてのそれぞれの計算された活性化スコアと、(ii)摂動シグネチャのセットにおける対応する化合物についてのそれぞれの摂動シグネチャのそれぞれの数値的活性化スコアとの間のそれぞれの差を使用して訓練されていないモデルを訓練することを更に含む。訓練は、差に回答して訓練されていないモデルと関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、化学化合物を目的の生理学的状態と関連付ける訓練されたモデルを得る。

10

#### 【0423】

本開示の更に別の態様は、1つ以上のプロセッサ、及び1つ以上のプロセッサによる実行のための1つ以上のプログラムを格納するメモリを有するコンピュータシステムを提供し、1つ以上のプログラムは、本明細書に開示される方法及び/又は実施形態のうちいずれかを実施するための命令を含む。いくつかの実施形態において、本開示の方法及び/又は実施形態のいずれかは、1つ以上のプロセッサ、及び1つ以上のプロセッサによって実行するための1つ以上のプログラムを格納するメモリを有するコンピュータシステムにおいて実施される。

20

#### 【0424】

本開示の別の態様は、コンピュータによって実行するように構成された1つ以上のプログラムを格納する非一時的なコンピュータ可読記憶媒体を提供し、1つ以上のプログラムは、本明細書に開示される方法のいずれかを実行するための命令を含む。

#### 【0425】

#### IV. 細胞構成要素モジュールの識別

いくつかの実施形態において、目的の生理学的状態と関連付けられた細胞構成要素モジュール132が識別される。このような方法は、図2及び図14と併せて本明細書で説明される。特に、図14Aのブロック1500を参照すると、いくつかの実施形態において、方法は、目的の生理学的状態と関連付けられた第1の細胞構成要素モジュール132を識別することを更に含む。

30

#### 【0426】

本開示のいくつかの実施形態に従って、細胞構成要素を目的の生理学的状態と関連付けるための方法200の例示的なワークフローは、図2A~図2Bを参照して提供される。

#### 【0427】

図2Aのブロック202及び図14Aのブロック1502を参照すると、方法は、1つ以上の第1のデータセットを電子形式で得ることを含む。図14Bのブロック1504を参照すると、1つ以上の第1のデータセットは、第1の複数の細胞におけるそれぞれの細胞の各々について、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含む。このようにして、複数のベクトルが得られる。

40

#### 【0428】

いくつかの実施形態において、目的の生理学的状態は、疾患であり、第1の複数の細胞が、複数の注釈付きの細胞状態によって立証されるように、疾患を代表する細胞、及び疾患を代表しない細胞を含む。

#### 【0429】

いくつかの実施形態において、図3Aのブロック300の目的の生理学的状態は、疾患と関連付けられた異常な細胞プロセスであり、第1の複数の細胞は、注釈付きの細胞状態

50

によって立証されるように、疾患を代表する細胞、及び疾患を代表しない細胞を含む。

【0430】

いくつかの実施形態において、図3Aのブロック300の目的の生理学的状態は、疾患と関連付けられた異常な細胞プロセスであり、第1の複数の細胞は、注釈付きの細胞状態によって立証されるように、疾患状態を代表する細胞、及び健康又は対照状態を代表する細胞を含む。

【0431】

いくつかの実施形態において、図3Aのブロック300の目的の生理学的状態は、複数の疾患と関連付けられた異常な細胞プロセスであり、第1の複数の細胞は、複数の注釈付きの細胞状態によって立証されるように、複数の細胞のサブセット、複数の疾患におけるそれぞれの疾患を代表する細胞のそれぞれのサブセットの各々を含む。

10

【0432】

図14Bのブロック1506を参照すると、いくつかの実施形態において、第1の複数の細胞は、2個、3個、4個、5個、6個、7個、8個、9個、10個、11個、12個、13個、14個、15個、16個、17個、18個、19個、20個、30個、40個、50個、60個、70個、80個、100個、200個、又は1000個以上の細胞を含み、複数（例えば、2個、3個、4個、5個、6個、7個、8個、9個、10個、11個、12個、13個、14個、15個、16個、17個、18個、19個、20個、30個、40個、50個、60個、70個、80個、100個、200個、又は1000個）の注釈付きの細胞状態を集合的に表す。

20

【0433】

図14Bのブロック1508を参照すると、いくつかの実施形態において、複数の細胞構成要素は、2個、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、50個、100個、500個、1000個、5000個、10,000個以上の細胞構成要素を含む。いくつかの実施形態において、複数の細胞構成要素は、2~10,000個又は細胞構成要素からなる。いくつかの実施形態において、複数の細胞構成要素は、100~10,000個又は細胞構成要素からなる。

【0434】

図2Aのブロック204を参照すると、方法は、複数のベクトルにアクセスすること、又はそれを形成することを含む。図14Aのブロック1510を参照すると、複数のベクトルにおけるそれぞれのベクトルの各々は、(i)複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii)対応する複数のエレメントを含む。図14Aのブロック1512を参照すると、対応する複数のエレメントにおけるそれぞれのエレメントの各々は、第1の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を表す対応するカウントを有する。

30

【0435】

ブロック206を参照すると、複数のベクトルは、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別するために使用される。複数の細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々は、複数の細胞構成要素のサブセットを含む。複数の細胞構成要素モジュールは、(i)複数の候補細胞構成要素モジュール、及び(ii)複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、複数の細胞構成要素モジュールは、10を超える細胞構成要素モジュールを含む。

40

【0436】

図14Bのブロック1514を参照すると、いくつかの実施形態において、複数の注釈付きの細胞状態における注釈付きの細胞状態は、曝露条件下（例えば、曝露期間、化合物の濃度、又は曝露期間及び化合物の濃度の組み合わせ）の化合物への第1の複数の細胞における細胞の曝露である。

【0437】

図14Bのブロック1518を参照すると、いくつかの実施形態において、複数の細胞

50

構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定の mRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである。

【0438】

図14Bのブロック1520を参照すると、いくつかの実施形態において、第1又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動 (FRET) 測定によって決定される。

【0439】

図14Bのブロック1522を参照すると、いくつかの実施形態において、第1又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、単一細胞リボ核酸 (RNA) 配列決定 (scRNA-seq)、scTag-seq、配列決定を使用したトランスポザーゼ-アクセス可能なクロマチンのための単一細胞アッセイ (scATAC-seq)、CyTOF/SCoP、E-MS/Abseq、miRNA-seq、CITE-seq、又はそれらの任意の組み合わせによって決定される。

10

【0440】

図14Bのブロック1524を参照すると、いくつかの実施形態、目的の生理学的状態は、疾患であり、第1の複数の細胞が、複数の注釈付きの細胞状態によって立証されるように、疾患を代表する細胞、及び疾患を代表しない細胞を含む。

20

【0441】

図14Bのブロック1526を参照すると、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別するために使用され、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々は、複数の細胞構成要素のサブセットを含む。複数の細胞構成要素モジュールは、(i) 複数の候補細胞構成要素モジュール、及び(ii) 複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、複数の細胞構成要素モジュールは、10を超える細胞構成要素モジュールを含む。

【0442】

図14Cのブロック1528を参照すると、いくつかの実施形態において、複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することは、複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、複数のベクトルに相関モデルを適用することを含む。いくつかの実施形態において、相関モデルは、グラフクラスタリングアルゴリズム (例えば、グラフクラスタリング方法は、ピアソン相関ベースの距離メトリック上のライデン (Leiden) クラスタリングであり、グラフクラスタリング方法は、ルーバン (Louvain) クラスタリングなどである)。

30

【0443】

図14Cのブロック1532を参照すると、いくつかの実施形態において、複数の細胞構成要素モジュールは、10~2000個、100~10000個、20~5000個、2~15,000個、80~5000個、100~500個の細胞構成要素モジュールからなる。いくつかの実施形態において、複数の細胞構成要素モジュールは、2~500個の細胞構成要素モジュールである。

40

【0444】

図14Cのブロック1534を参照すると、いくつかの実施形態において、複数の細胞構成要素は、10~2000個、100~10000個、20~5000個、2~15,000個、80~5000個、100~500個の細胞構成要素からなる。いくつかの実施形態において、複数の細胞構成要素は、2~500個の細胞構成要素である。

【0445】

図14Cのブロック1536を参照すると、いくつかの実施形態において、複数の構成

50



要素モジュールにおける候補細胞構成要素モジュールの各々は、200～300個の細胞構成要素からなる。

【0446】

図2Aのブロック208及び図14Cのブロック1538を参照すると、方法は、1つ以上の第2のデータセットを電子形式で得ることを含む。1つ以上の第2のデータセットは、第2の複数の細胞におけるそれぞれの細胞の各々について、第2の複数の細胞が、20個以上の細胞を含み、目的の生理学的状態を通知する複数の共変量を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含むか、又は集合的に含む。したがって、細胞構成要素カウントデータ構造が得られ、細胞構成要素カウントデータ構造は、(i) 第2の複数の細胞、及び(ii) 複数の細胞構成要素又はその表現によって次元決定される。

10

【0447】

図14Cのブロック1540を参照すると、いくつかの実施形態において、複数の共変量は、細胞バッチ、細胞ドナー、細胞型、疾患状態、又は化学化合物への曝露を含む。

【0448】

図2Bのブロック210及びブロック1542 図14Dを参照すると、活性化データ構造は、複数の細胞構成要素又はその表現を共通次元として使用して、細胞構成要素カウントデータ構造及び潜在表現を組み合わせることによって形成される。活性化データ構造は、複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、第2の複数の細胞における細胞の各々について、それぞれの活性化重みを含む。

20

【0449】

図2Bのブロック212及び図14Dのブロック1544を参照すると、方法は、(i) 活性化データ構造を候補モデルに入力したときに、活性化データ構造内に表される細胞構成要素モジュールの各々における複数の共変量における各共変量の不在又は存在の予測と、(ii) 細胞構成要素モジュールの各々における各共変量の実際の不在又は存在との間の差を使用して候補細胞構成要素モデルを訓練することを更に含む。訓練することは、差に応答して候補細胞構成要素モデルと関連付けられた複数の共変量重みを調整し、複数の共変量重みは、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの共変量の各々について、それぞれの共変量が、活性化データ構造にわたって、それぞれの細胞構成要素モジュールと相関するかどうかを示す対応する重みを含む。

30

【0450】

図14Dのブロック1546を参照すると、候補細胞構成要素モデルを訓練することは、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、複数のコスト関数におけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する。

【0451】

したがって、図2Cのブロック214及び図14Dのブロック1548を参照すると、複数の共変量重みは、候補細胞構成要素モデルを訓練する際に、複数の候補細胞構成要素モジュールにおける第1の細胞構成要素モジュールを識別するために使用され、複数の候補細胞構成要素モジュールにおける第1の細胞構成要素モジュールは、目的の生理学的状態と関連付けられる。

40

【0452】

いくつかの実施形態において、第1及び/又は第2の複数の細胞は、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも30個、少なくとも40個、少なくとも50個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも1000個、少なくとも2000個、少なくとも3000個、少なくとも4000個、少なくとも5000個、少なくとも10,000個、少なくとも20,000個、少なくとも30,000

50

0個、少なくとも50,000個、少なくとも80,000個、少なくとも100,000個、少なくとも500,000個、又は少なくとも100万個の細胞を含む。いくつかの実施形態において、第1及び/又は第2の複数の細胞は、500万個以下、100万個以下、500,000個以下、100,000個以下、50,000個以下、10,000個以下、5000個以下、1000個以下、500個以下、200個以下、100個以下、又は50個以下の細胞を含む。いくつかの実施形態において、第1及び/又は第2の複数の細胞は、5~100個、10~50個、20~500個、200~10,000個、1000~100,000個、50,000~500,000個、又は10,000~100万個の細胞を含む。いくつかの実施形態において、第1及び/又は第2の複数の細胞は、5個以上の細胞から始まり、500万個以下の細胞で終わる別の範囲内にある。

10

## 【0453】

いくつかの実施形態において、第2の複数の細胞は、第1の複数の細胞に含まれる細胞を含まない。いくつかの実施形態において、第2の複数の細胞は、第1の複数の細胞に含まれる細胞の一部又は全てを含む。

## 【0454】

いくつかの実施形態において、複数の注釈付きの細胞状態は、細胞表現型、細胞挙動、疾患状態、遺伝子変異、遺伝子若しくは遺伝子産物の摂動(例えば、ノックダウン、サイレンシング、過剰発現など)、及び/又は化合物への曝露のうちの一つ以上を含む。いくつかの実施形態において、複数の注釈付きの細胞状態における注釈付きの細胞状態は、曝露条件下での化合物への第1の複数の細胞における細胞の曝露である。例えば、細胞の曝露は、一つ以上の化合物での細胞の任意の処理を含む。いくつかの実施形態において、一つ以上の化合物は、例えば、小分子、生物製剤、治療剤、タンパク質、小分子と組み合わされたタンパク質、ADC、核酸(例えば、siRNA、干渉RNA、cDNA過剰発現野生型及び/若しくは変異体shRNA、cDNA過剰発現野生型及び/若しくは変異体ガイドRNA(例えば、Cas9系若しくは他の細胞成分編集系)など)、並びに/又は前述のいずれかの任意の組み合わせを含む。いくつかの実施形態において、曝露条件は、曝露期間、化合物の濃度、又は曝露期間及び化合物の濃度の組み合わせである。いくつかの実施形態において、化合物は、上記の「化合物」と題されるセクションにおいてなど、本明細書に記載される実施形態のいずれでもある。

20

## 【0455】

いくつかの実施形態において、複数の注釈付きの細胞状態は、細胞バッチ、細胞ドナー、細胞型、細胞株、疾患状態、時点、複製、及び/又は関連するメタデータの一つ以上の兆候を含む。いくつかの実施形態において、複数の注釈付きの細胞状態は、実験データ(例えば、フローサイトメトリーの読み出し、イメージング及び顕微鏡注釈、細胞構成要素データなど)を含む。いくつかの実施形態において、複数の注釈付きの細胞状態は、一つ以上の遺伝子マーカー(例えば、コピー数バリエーション、単一ヌクレオチドポリモーフィズム、多ヌクレオチド多型、挿入、欠失、遺伝子融合、マイクロサテライト不安定性状態、増幅、及び/又はアイトフォーム)を含む。いくつかの実施形態において、複数の注釈付きの細胞状態は、本明細書に開示される共変量のうちの一つ以上及び/又は本明細書に開示される目的の生理学的状態のうちの一つ以上、例えば、上記の「生理学的状態」と題されるセクションなどを含む。

30

40

## 【0456】

本明細書に開示される任意の細胞構成要素及び/又は任意の細胞構成要素モジュール、並びにそれらの任意の実施形態、置換、修飾、追加、欠失、及び/又は組み合わせは、上記の「細胞構成要素及び細胞構成要素モジュール」と題されたセクションに記載されるように、細胞構成要素モジュールの識別のために企図される。例えば、いくつかの実施形態において、複数の細胞構成要素における細胞構成要素の各々は、特定の遺伝子、遺伝子に関連する特定のmRNA、炭水化物、脂質、エピジェネティック特徴、代謝産物、タンパク質、又はそれらの組み合わせである。いくつかの実施形態において、複数の細胞構成要素は、100~8,000個の細胞構成要素からなる。いくつかの実施形態において、複

50

数の細胞構成要素モジュールは、10～2000個の細胞構成要素モジュールからなる。いくつかの実施形態において、複数の構成要素モジュールにおける候補細胞構成要素モジュールの各々は、200～300個の細胞構成要素からなる。

【0457】

いくつかの実施形態において、それぞれの細胞構成要素の対応する存在量は、上記に開示される細胞構成要素のいずれかの存在量を含む。

【0458】

いくつかの存在量カウント技術（例えば、細胞構成要素測定技術）のうちのいずれか1つを使用して、それぞれの細胞の各々におけるそれぞれの細胞構成要素の各々についての対応する存在量を得ることができる。例えば、表1は、本開示のいくつかの実施形態に従う、単一細胞の細胞構成要素測定のための非限定的な技術を列挙する。

10

【0459】

いくつかの実施形態において、それぞれの細胞構成要素の対応する存在量は、蛍光、化学発光、電気シグナル検出、ポリメラーゼ連鎖反応（PCR）、逆転写酵素ポリメラーゼ連鎖反応（RT-PCR）、デジタル液滴PCR（ddPCR）、固体状態ナノポア検出、RNAスイッチ活性化、ノーザンブロット、及び/又は遺伝子発現の連続分析（SAGE）を介したマイクロアレイ分析を含む1つ以上の方法を使用して決定される。いくつかの実施形態において、第1又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、比色測定、蛍光測定、発光測定、又は共鳴エネルギー移動（FRET）測定によって決定される。

20

【0460】

いくつかの実施形態において、第1及び/又は第2の複数の細胞におけるそれぞれの細胞における遺伝子発現は、細胞を配列決定し、次いで配列決定中に識別された各遺伝子転写物の量をカウントすることによって測定することができる。いくつかの実施形態において、配列決定及び定量化された遺伝子転写産物は、mRNAなどのRNAを含む。いくつかの実施形態において、配列決定及び定量化された遺伝子転写産物は、タンパク質（例えば、転写因子）などのmRNAの下流産物を含む。一般に、本明細書で使用される場合、「遺伝子転写産物」という用語は、翻訳後修飾を含む、遺伝子転写又は翻訳の任意の下流産物を示すために使用されてもよく、「遺伝子発現」は、一般に、遺伝子転写産物の任意の尺度を指すために使用されてもよい。

30

【0461】

いくつかの実施形態において、それぞれの細胞構成要素の対応する存在量はRNA存在量（例えば、遺伝子発現）であり、それぞれの細胞構成要素の存在量は、それぞれの遺伝子に対応する1つ以上の核酸分子のポリヌクレオチドレベルを測定することによって決定される。それぞれの遺伝子の転写レベルは、第1及び/又は第2の複数の細胞におけるそれぞれの細胞中に存在するmRNA又はそれに由来するポリヌクレオチドの量から決定することができる。ポリヌクレオチドは、マイクロアレイ分析、ポリメラーゼ連鎖反応（PCR）、逆転写酵素ポリメラーゼ連鎖反応（RT-PCR）、ノーザンブロット、遺伝子発現の連続分析（SAGE）、RNAスイッチ、RNAフィンガープリンティング、リガーゼ連鎖反応、Qベータレプリカーゼ、等温増幅法、鎖置換増幅、転写ベース増幅システム、ヌクレアーゼ保護アッセイ（Siヌクレアーゼ又はRNAse保護アッセイ）、及び/又は固体状態ナノポア検出を含むが、これらに限定されない、様々な方法によって検出及び定量することができる。例えば、Draghici, Data Analysis Tools for DNA Microarrays, Chapman and Hall/CRC, 2003、Simon et al., Design and Analysis of DNA Microarray Investigations, Springer, 2004、Real-Time PCR: Current Technology and Applications, Logan, Edwards, and Saunders eds., Caister Academic Press, 2009、Bustin A-Z of Quantitative PCR (IUL Bio

40

50

technology, No. 5), International University Line, 2004、Velculescu et al., (1995) Science 270: 484 - 487、Matsumura et al., (2005) Cell. Microbiol. 7: 11 - 18、Serial Analysis of Gene Expression (SAGE): Methods and Protocols (Methods in Molecular Biology), Humana Press, 2008を参照されたく、これらの各々は参照によりその全体が本明細書に組み込まれる。

#### 【0462】

いくつかの実施形態において、それぞれの細胞構成要素の対応する存在量は、発現RNA又はそれに由来する核酸(例えば、RNAポリメラーゼプロモーターを組み込んだcDNAに由来するcDNA又は増幅RNA)から、天然核酸分子、及び合成核酸分子を含む、第1及び/又は第2の複数の細胞におけるそれぞれの細胞から得られる。したがって、いくつかの実施形態において、それぞれの細胞構成要素の対応する存在量は、総細胞RNA、ポリ(A)+メッセンジャーRNA(mRNA)若しくはその画分、細胞質mRNA、又はcDNAから転写されたRNA(例えば、cRNA)などの非限定的な供給源から得られる。総RNA及びポリ(A)+RNAを調製するための方法は、当該技術分野で周知であり、一般に、例えば、Sambrook, et al., Molecular Cloning: A Laboratory Manual (3rd Edition, 2001)に記載されている。RNAは、グアニジンチオシアナート溶解後のCsCl遠心分離(例えば、Chirgwin et al., 1979, Biochemistry 18: 5294 - 5299を参照されたい)、シリカゲルベースのカラム(例えば、RNeasy (Qiagen, Valencia, Calif.)若しくはStrataPrep (Stratagene, La Jolla, Calif.))を使用して、又はAusubel et al., eds., 1989, Current Protocols In Molecular Biology, Vol. III, Green Publishing Associates, Inc., John Wiley & Sons, Inc., New York, pp. 13.12.1 - 13.12.5)に記載されているフェノール及びクロロホルムを使用して目的の細胞から抽出することができる。ポリ(A)+RNAは、例えば、オリゴ-dTセルロースを用いた選択によって、又は代替的に、全細胞RNAのオリゴ-dTプライミング逆転写によって選択することができる。RNAは、当該技術分野で既知の方法によって、例えば、ZnCl<sub>2</sub>とのインキュベーションによって断片化して、RNAの断片を生成することができる。

#### 【0463】

いくつかの実施形態において、第1及び/又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、配列決定によって決定される。いくつかの実施形態において、第1及び/又は第2の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、単一細胞リボ核酸(RNA)配列決定(scRNA-seq)、scTag-seq、配列決定を使用したトランスクリプターゼ-アクセス可能なクロマチンのための単一細胞アッセイ(scATAC-seq)、CyTOF/SCoP、E-MS/Abseq、miRNA-seq、CITE-seq、及びそれらの任意の組み合わせによって決定される。

#### 【0464】

細胞構成要素存在量測定技術は、測定される所望の細胞構成要素に基づいて選択することができる。例えば、scRNA-seq、scTag-seq、及びmiRNA-seqを使用して、RNA発現を測定することができる。具体的には、scRNA-seqはRNA転写産物の発現を測定し、scTag-seqは希少なmRNA種の検出を可能にし、miRNA-seqはマイクロRNAの発現を測定する。CyTOF/SCoP及びE-MS/Abseqを使用して、細胞内のタンパク質発現を測定することができる。CITE-seqは、細胞における遺伝子発現及びタンパク質発現の両方を同時に測定し、

s c A T A C - s e q は、細胞におけるクロマチンコンフォメーションを測定する。以下の表 1 は、上記の細胞構成要素存在量測定技術の各々を実施するための例示的なプロトコルを提供する。

【 0 4 6 5 】

【 表 1 】

表 1 - 測定プロトコルの例

技術	プロトコル
RNA-seq	Olsen <i>et al.</i> , (2018), "Introduction to Single-Cell RNA Sequencing," <i>Current protocols in molecular biology</i> 122(1), pg.57.
Tag-seq	Rozenberg <i>et al.</i> , (2016), "Digital gene expression analysis with sample multiplexing and PCR duplicate detection: A straightforward protocol," <i>BioTechniques</i> , 61(1), pg.26.
ATAC-seq	Buenrostro <i>et al.</i> , (2015), "ATAC-seq: a method for assaying chromatic accessibility genome-wide," <i>Current protocols in molecular biology</i> , 109(1), pg.21.
miRNA-seq	Faridani <i>et al.</i> , (2016), "Single-cell sequencing of the small-RNA transcriptome," <i>Nature biotechnology</i> , 34(12), pg.1264.
CytoTOF/SCoPE-MS/Abseq	Bandura <i>et al.</i> , (2009), "Mass cytometry: technique for real time single cell multitarget immunoassay based on inductively coupled plasma time-of-flight mass spectrometry," <i>Analytic chemistry</i> , 81(16), pg.6813.  Budnik <i>et al.</i> , (2018), "SCoPE-ME: mass spectrometry of single mammalian cells quantifies proteome heterogeneity during cell differentiation," <i>Genome biology</i> , 19(1), pg.161.  Shahi <i>et al.</i> , (2017), "Abseq: Ultrahigh-throughput single cell protein profiling with droplet microfluidic barcoding," <i>Scientific reports</i> , 7, pg.44447.
CITE-seq	Stoeckius <i>et al.</i> , (2017), "Simultaneous epitope and transcriptome measurement in single cells," <i>Nature Methods</i> , 14(9), pg.856.

10

20

30

【 0 4 6 6 】

いくつかの実施形態において、複数の細胞構成要素は、単一の時点で測定される。いくつかの実施形態において、複数の細胞構成要素は、複数の時点で測定される。例えば、いくつかの実施形態において、複数の細胞構成要素は、細胞状態遷移（例えば、分化プロセス、化合物への曝露に対する応答、発生プロセスなど）全体にわたる複数の時点で測定される。

【 0 4 6 7 】

本開示は、細胞（例えば、単一細胞）から得られる他の細胞構成要素の測定値を使用する類似の方法を包含するため、これは例示であり、限定ではないことを理解されたい。本開示は、本開示に記載される方法を実施する個人又は組織によって実施される実験作業から直接得られた測定値を使用する方法、並びに例えば、他者によって実施される実験作業の結果の報告から間接的に得られ、第三者の出版物、データベース、請負業者によって実施されるアッセイ、又は開示される方法を実施するのに有用な好適な入力データの他の供給源で報告されたデータを含む、任意の手段又は機構を通じて利用可能にされた測定値を使用する方法を包含することを更に理解されたい。

40

【 0 4 6 8 】

いくつかの実施形態において、第 1 及び / 又は第 2 の複数の細胞（例えば、1 つ以上の複数の第 1 のデータセット及び / 又は 1 つ以上の第 2 のデータセット）における複数の細

50

胞構成要素に対する対応する存在量は、前処理される。いくつかの実施形態において、前処理は、フィルタリング、正規化、マッピング（例えば、参照配列に対する）、定量化、スケールリング、デコンボリューション、クリーニング、次元縮小、変換、統計分析、及び/又は集約のうちの一つ以上を含む。

【0469】

例えば、いくつかの実施形態において、複数の細胞構成要素は、所望の品質、例えば、核酸配列のサイズ及び/若しくは品質、又はそれぞれの細胞構成要素についての最小及び/若しくは最大存在量値に基づいてフィルタリングされる。いくつかの実施形態において、フィルタリングは、Skewerなどの様々なソフトウェアツールによって一部又はその全体が実施される。Jiang, H. et al., BMC Bioinformatics 15(182): 1-12 (2014)を参照されたい。いくつかの実施形態において、複数の細胞構成要素は、例えば、AfterQC、Kraken、RNA-SeqQC、FastQC、又は別の同様のソフトウェアプログラムなどの配列決定データQCソフトウェアを使用して、品質管理のためにフィルタリングされる。いくつかの実施形態において、複数の細胞構成要素は、例えば、ブルダウン、増幅、及び/又は配列決定バイアス（例えば、マップビリティ、GCバイアスなど）を考慮するために正規化される。例えば、Schwartz et al., PLoS ONE 6(1): e16685 (2011) and Benjamini and Speed, Nucleic Acids Research 40(10): e72 (2012)を参照されたく、その内容は全ての目的のために参照によりその全体が本明細書に組み込まれる。いくつかの実施形態において、前処理は、細胞構成要素のサブセットを複数の細胞構成要素から除去する。いくつかの実施形態において、複数の細胞構成要素について対応する存在量を前処理することは、高い信号対ノイズ比を改善する（例えば、低下させる）。

【0470】

いくつかの実施形態において、前処理は、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量と参照の存在量との比較を実施することを含む。いくつかの実施形態において、参照存在量は、例えば、正常試料、適合した試料、参照存在量値を含む参照データセット、ハウスキーピング遺伝子などの参照細胞構成要素、及び/又は参照標準から得られる。いくつかの実施形態において、細胞構成要素存在量のこの比較は、平均検定、ウィルコクソンランクサム検定（マンホイットニーU検定）、t検定、ロジスティック回帰、及び一般化線形モデルの差異を含むが、これらに限定されない任意の差次的発現試験を使用して実施される。当業者は、細胞構成要素存在量の比較及び/又は正規化のために他のメトリックも可能であることを理解するであろう。

【0471】

したがって、いくつかの実施形態において、一つ以上の第1のデータセット及び/又は一つ以上の第2のデータセットにおけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量は、限定されないが、生の存在量値、絶対存在量値（例えば、転写物数）、相対的な存在量値（例えば、相対的な蛍光単位、トランスクリプトーム分析、及び/若しくは遺伝子セット発現分析（GSEA））、化合物若しくは集合的な存在量値、変換された存在量値（例えば、 $\log_2$ 及び/若しくは $\log_{10}$ 変換）、参照（例えば、通常の試料、適合した試料、参照データセット、ハウスキーピング遺伝子、及び/若しくは参照標準）に対する変化（例えば、倍数若しくは $\log$ 変化）、標準化された存在量値、中心傾向の尺度（例えば、平均、中央値、モード、加重平均、加重中央値、及び/若しくは加重モード）、分散の尺度（例えば、不一致、標準偏差、及び/若しくは標準誤差）、調整された存在量値（例えば、正規化された、スケールリングされた、及び/若しくは誤差訂正された）、次元低減された存在量値（例えば、主成分ベクトル及び/若しくは潜在成分）、並びに/又はそれらの組み合わせを含む、様々な形態のうちの一つを含む。次元縮小技術を使用して細胞構成要素の存在量を得るための方法は、主成分分析、因子分析、線形判別分析、多次元スケールリング、等角特徴マッピング、局所線形埋め込み、ヘシアン固有マッピング、スペクトル埋め込み、t分布確率論的隣接埋め込み、並びに/又

10

20

30

40

50

はそれらの任意の置換、追加、欠失、修飾、及び/若しくは組み合わせを含むが、これらに限定されない、当該技術分野において既知であり、以下で更に詳細に説明されることは、当業者に明らかであろう。例えば、参照によりその全体が本明細書に組み込まれる、S umithra et al., 2015, "A Review of Various Linear and Non Linear Dimensionality Reduction Techniques," Int J Comp Sci and Inf Tech, 6 (3), 2354 - 2360を参照されたい。

【0472】

いくつかの実施形態において、複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別することは、複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して、複数のベクトルに相関モデルを適用することを含む。

10

【0473】

いくつかの実施形態において、相関モデルは、クラスタリング方法（例えば、クラスタリングモデル）を含む。いくつかの実施形態において、相関モデルは、グラフクラスタリング方法（例えば、モデル）及び/又は非グラフクラスタリング方法を含む。いくつかの実施形態において、グラフクラスタリング方法は、ピアソン相関ベースの距離メトリック上のライデン（Leiden）クラスタリングである。いくつかの実施形態において、グラフクラスタリング方法は、ルーバン（Louvain）クラスタリングである。

【0474】

例えば、いくつかの実施態様では、方法は、相関ベースのコスト関数の適用を含む。相関ベースのコスト関数を最適化することは、細胞構成要素（例えば、遺伝子）間の最近傍関係を定義する最近傍グラフを計算することと、各細胞構成要素を、各細胞内の細胞構成要素についての存在量カウント（例えば、発現値）を格納することによって形成されるベクトルによって表すことと、細胞構成要素間の相関を計算することと、を含む。互いに高い相関を有する細胞構成要素は、最近傍であると判定され、グラフクラスタリング方法（例えば、ライデン（Leiden）及び/又はルーバン（Louvain））を使用してグラフをクラスタリングすることによって、細胞構成要素モジュールを形成するために使用される。

20

【0475】

いくつかのクラスタリング技術のうちの任意の1つを使用することができ、その例としては、階層的クラスタリング、k平均クラスタリング、及び密度ベースのクラスタリングが含まれるが、これらに限定されない。一実施形態において、階層的密度ベースのクラスタリングが使用される（HDBSCANと称され、例えば、Campello et al., (2015). Hierarchical density estimates for data clustering, visualization, and outlier detection. ACM Trans Knowl Disc Data, 10 (1), 5を参照されたい）。別の実施形態において、ルーバン（Louvain）クラスタリングなどのコミュニティ検出に基づくクラスタリングが使用される（例えば、Blondel et al., (2008). Fast unfolding of communities in large networks. J stat mech: theor exp, 2008 (10), P10008を参照されたい）。なお別の実施形態において、ライデン（Leiden）クラスタリングが使用される。ライデン（Leiden）アルゴリズムは、個々のノードをコミュニティ間で移動してパーティションを決定し、パーティションを洗練し、洗練されたパーティションに基づいて集約ネットワークを作成することによって進行する。集約ネットワークは、プロセスの以前のステップで決定された未洗練のパーティションに基づいて更に分割され、新しいパーティションは、各集約ネットワーク内の個々のノードを移動することによって洗練される。例えば、Traag et al., (2019), "From Louvain to Leiden: guaranteeing well-connected communi

30

40

50

ties, " Sci Rep 9 : 5 2 3 3 , doi : 1 0 . 1 0 3 8 / s 4 1 5 9 8 - 0 1 9 - 4 1 6 9 5 - z を参照されたい。なお別の実施形態において、拡散経路アルゴリズムが使用される。

【 0 4 7 6 】

一般に、ルーバン ( Louvain ) クラスタリング及び / 又はライデン ( Leiden ) クラスタリングなどのクラスタリングは、ハードパーティショニング技術を使用し、各エレメント ( 例えば、各細胞構成要素 ) は、重複することなく単一のクラスタに一意に割り当てられる。しかしながら、任意の1つの特定の理論に拘束されることなく、細胞プロセス ( 例えば、目的の生理学的状態と関連付けられる ) は、細胞内の細胞構成要素のネットワーク間の複雑かつ動的な相互作用によって特徴付けられ得、例えば、単一の遺伝子は、細胞内の2つ、3つ、4つ、又はそれ以上の細胞プロセスにおいて、任意の数の同じ又は異なるプロセス及び経路において同様に機能する任意の数の他の遺伝子と組み合わせることで役割を果たすことができる。したがって、細胞内活性の複雑さと並行して、第1のモジュールへの細胞構成要素のクラスタリングは、必ずしも他のモジュールを除外する必要はない。したがって、いくつかの実施形態において、細胞構成要素モジュールの識別は、細胞構成要素の重複するサブセットを有するモジュールを得ることを含む。

10

【 0 4 7 7 】

相関ベースのモデルを使用してハードパーティショニング技術を利用することの代わりに、又はそれに加えて、いくつかの実施形態において、複数のベクトルを使用して複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々を識別することは、複数の細胞構成要素の表現を複数の次元縮小構成要素として生成する辞書学習モデルを含む。いくつかの実施形態において、辞書学習モデルは、L0正規化オートエンコーダである。これらのモデルの利点は、モジュールと細胞構成要素との間に1:1の対応を強制しないが、細胞構成要素が同時にいくつかのモジュールに現れることを可能にすることである。

20

【 0 4 7 8 】

例えば、いくつかの実施形態において、方法は、スペアオートエンコーダコスト関数の適用を含む。いくつかのそのような例では、スパースオートエンコーダのコスト関数を最適化することは、pytorch又はtensorflowに実装されているような標準訓練を使用して、その重みのL0正規化、及び再構築損失を伴う1層オートエンコーダを訓練することを含む。

30

【 0 4 7 9 】

限定されないが、ファジーK平均、重複K平均 ( OKM )、重み付けOKM ( WOKM )、重複分割クラスタ ( OPC )、及びマルチクラスタ重複K平均拡張 ( MCKE )、並びに / 又はそれらの任意の変形若しくは組み合わせを含む、重複分割アルゴリズムの他の方法が可能である。

【 0 4 8 0 】

いくつかの実施形態において、統計技術は、1つ以上の第1のデータセットに符号化された潜在情報の形状を維持しながら、高次元データ ( 例えば、複数の注釈付きの細胞状態を集散的に表す第1の複数の細胞における細胞の各々について、複数の細胞構成要素モジュールにわたる複数の細胞構成要素の存在量 ) を低次元空間に圧縮するために使用され得る。例えば、図4の上部パネルに示されるように、カウントマトリックスは、第1の複数の細胞における細胞の各々について、複数の細胞構成要素における細胞構成要素の各々について、対応するカウント ( 例えば、存在量 ) を含む。カウントマトリックスは、異なる注釈付きの細胞状態 ( 例えば、細胞型、曝露条件、疾患など ) の条件下でのそれらの対応する存在量の類似性に基づいて、第1の複数の細胞にわたる細胞構成要素のクラスタリングを表す低次元空間にデータが縮小される、図4の下部パネルに示される潜在表現に変換することができる。したがって、クラスタ化された細胞構成要素は、細胞構成要素モジュールとして表され、潜在表現では、複数の細胞状態にわたる挙動の類似性を符号化する。

40

【 0 4 8 1 】

図4に示される潜在表現を再び参照すると、各行 - 列グループ化におけるエントリー内

50



の値は、元の入力データセットに基づいて次元数の減少によって決定される。例えば、各エントリーは、それぞれの列によって表されるそれぞれの細胞構成要素の各々について、それぞれの行によって表されるそれぞれの細胞構成要素モジュールに含まれる複数の細胞構成要素のサブセット（例えば、重み $w_{1-1}$ 、重み $w_{1-2}$ など）におけるメンバーシップの表示を含むことができる。特に、いくつかの実施形態において、各エントリーは、それぞれの細胞構成要素がそれぞれのモジュールに含まれるかどうかを示す重みである。いくつかの実施形態において、重みは、メンバーシップの2値表示である（例えば、それぞれのモジュールにおける存在又は不在は、それぞれ1又は0で示される）。いくつかの実施形態において、重みは、それぞれのモジュールに対する細胞構成要素の相対的な重要性（例えば、メンバーシップの確率及び/又は相関）を示すようにスケールされる。

10

## 【0482】

いくつかの実施形態において、潜在表現におけるそれぞれの次元は、それぞれの細胞構成要素の表現に対応する。細胞構成要素の表現は、例えば、潜在表現マトリックス内のそれぞれのエントリー（例えば、重み）が複数の細胞構成要素に対応する場合などの、細胞構成要素の非線形表現から生じ得る。細胞構成要素の表現を含む他の実施形態は、主成分分析を使用して得られた潜在表現を含み、各主成分は、複数の細胞構成要素に対応するデータの分散及び/又は他の変換を表す。

## 【0483】

いくつかの実施形態において、次元数削減技術は、データのいくつかの非可逆圧縮をもたらす。しかしながら、結果として生じる潜在表現（例えば、潜在表現118）は、計算記憶サイズにおいてより小さく、したがって、モデル訓練などの他の下流技術と併せて分析するためのより少ない計算処理能力を必要とする。したがって、潜在表現における複数の細胞構成要素モジュールの配置は、現代のコンピューティングデバイスを使用して、現在開示されている方法の計算実現可能性を増加させる。

20

## 【0484】

様々な次元数削減技術を使用することができる。いくつかの実施形態において、次元数削減は、主成分（PCA）、ランダム投影、独立成分分析、特徴選択、因子分析、Sammonマッピング、曲線成分分析、確率的隣接埋め込み（SNE）、アイソマップ、最大分散展開、局所線形埋め込み、t-SNE、非負のマトリックス因子分解、カーネル主成分分析、グラフベースのカーネル主成分分析、線形判別分析（LDA）、一般化判別分析、一様多様体近似及び投影（UMAP）、LargeVis、Laplacian Eigenmap、拡散マップ、ネットワーク（例えば、ニューラルネットワーク）技術、及び/又はフィッシャーの線形判別分析である。例えば、Fodor, 2002, "A survey of dimension reduction techniques," Center for Applied Scientific Computing, Lawrence Livermore National, Technical Report UCR L-ID-148494, Cunningham, 2007, "Dimension Reduction," University College Dublin, Technical Report UCD-CSI-2007-7, Zahorian et al., 2011, "Nonlinear Dimensionality Reduction Methods for Use with Automatic Speech Recognition," Speech Technologies. doi:10.5772/16863. ISBN 978-953-307-996-7、及びLakshmi et al., 2016, "2016 IEEE 6th International Conference on Advanced Computing (IACC)," pp. 31-34. doi:10.1109/IACC.2016.16, ISBN 978-1-4673-8286-1を参照されたく、それらの各々は参照により本明細書に組み込まれる。したがって、いくつかの実施形態において、次元数削減は、主成分分析（PCA）であり、それぞれの抽出された次元数削減成分の各々は、PCAによって導出されたそれぞれの主成分を含む。そのような実施形態にお

30

40

50

いて、複数の主成分における主成分の数は、PCAによって計算された主成分の閾値数に制限され得る。主成分の閾値数は、例えば、少なくとも5、少なくとも10、少なくとも20、少なくとも50、少なくとも100、少なくとも1000、少なくとも1500、又は任意の他の数であり得る。いくつかの実施形態において、PCAによって計算された各主成分は、PCAによって固有値を割り当てられ、第1の複数の抽出された特徴の対応するサブセットは、最も高い固有値を割り当てられた主成分の閾値数に限定される。複数の細胞構成要素ベクトルにおけるそれぞれの細胞構成要素ベクトルの各々について、複数の次元削減コンポーネントをそれぞれの細胞構成要素ベクトルに適用して、複数の次元削減コンポーネントにおけるそれぞれの次元削減コンポーネントの各々についての次元削減コンポーネント値を含む対応する次元削減ベクトルを形成する。これは、複数の細胞構成要素ベクトルから、対応する複数の次元削減ベクトルを形成し、それによって、潜在表現に配置された複数の細胞構成要素モジュールを形成する。

10

#### 【0485】

いくつかの実施形態において、方法は、潜在表現で配置された複数の細胞構成要素モジュールを使用して多様体学習を実施することを更に含む。一般に、多様体学習は、データセットにおける最大変動を決定することによって、高次元データの低次元構造を説明するために使用される。例としては、限定されないが、力指向レイアウト (Fruchterman, T. M., & Reingold, E. M. (1991). Graph drawing by force-directed placement. *Software: Practice and experience*, 21(11), 1129-1164) (例えば、Force Atlas 2)、t分布型確率的近傍埋め込み法 (t-SNE)、局所線形埋め込み (Roweis, S. T., & Saul, L. K. (2000). Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500), 2323-2326)、局所線形アイソメトリックマッピング (ISOMAP, Tenenbaum, J. B., De Silva, V., & Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500), 2319-2323)、カーネルPCA、グラフベースのカーネルPCA、親和性ベースの軌道埋め込みのための熱拡散の可能性 (Potential of Heat-Diffusion for Affinity Based Trajectory Embedding) (PHATE)、一般化判別分析 (GDA)、一様多様体近似及び投影 (UMAP)、又はカーネル判別分析が挙げられるが、これらに限定されない。判別分析は、特に、各細胞の特定の細胞型についていくつかの情報が事前に知られている場合に使用され得る。力指向レイアウトは、基礎となる細胞プロセスから生じる基礎となるデータの非線形の態様を符号化する新しい低次元を識別する能力のために、様々な特定の実施形態において有用である。力指向レイアウトは、データを最もよく表す縮小された次元を決定するための機構として、物理ベースのモデルを使用する。例として、力指向レイアウトは、本実施形態において、1つ以上の第1のデータセットにおける各細胞に「反発」力が割り当てられ、第1の複数の細胞にわたって計算されるときに、これらの競合する「力」の下で一緒に「拡散」するデータのセクタを識別するグローバルな「重力」が存在する物理シミュレーションの形態を使用する。力指向レイアウトは、データの構造についての仮定をほとんど行わず、ノイズ除去アプローチを課さない。

20

30

40

#### 【0486】

多様体学習は、例えば、Wang et al., 2004, "Adaptive Manifold Learning," *Advances in Neural Information Processing Systems* 17に更に記載され、これはその全体が参照により本明細書に組み込まれる。

#### 【0487】

50

いくつかの実施形態において、複数の共変量は、細胞バッチ、細胞ドナー、細胞型、疾患状態、又は化学化合物への曝露を含む。いくつかの実施形態において、複数の共変量は、第2の複数の細胞における1つ以上の細胞に関連する時点、複製、及び/又は関連するメタデータの1つ以上の表示を含む。いくつかの実施形態において、複数の共変量は、実験データ（例えば、フローサイトメトリーの読み出し、イメージング及び顕微鏡注釈、細胞構成要素データなど）を含む。いくつかの実施形態において、複数の共変量は、第2の複数の細胞における1つ以上の細胞に特徴的な1つ以上の遺伝子マーカー（例えば、コピー数バリエーション、単一ヌクレオチドポリモーフィズム、多ヌクレオチド多型、挿入、欠失、遺伝子融合、マイクロサテライト不安定性状態、増幅、及び/又はアイソフォーム）を含む。いくつかの実施形態において、複数の共変量は、第2の複数の細胞における1つ以上の細胞についての細胞表現型、細胞挙動、疾患状態、遺伝子変異、遺伝子若しくは遺伝子産物の摂動（例えば、ノックダウン、サイレンシング、過剰発現など）、及び/又は曝露条件のうちの1つ以上を含む。

10

## 【0488】

例えば、いくつかの実施形態において、共変量は、曝露条件下での化合物への第2の複数の細胞における細胞の曝露又は曝露に対する応答である。いくつかの実施形態において、細胞の曝露は、1つ以上の化合物での細胞の任意の処理を含む。いくつかの実施形態において、1つ以上の化合物は、例えば、小分子、生物製剤、治療剤、タンパク質、小分子と組み合わされたタンパク質、ADC、核酸（例えば、siRNA、干渉RNA、cDNA過剰発現野生型及び/若しくは変異体shRNA、cDNA過剰発現野生型及び/若しくは変異体ガイドRNA（例えば、Cas9系若しくは他の細胞成分編集系）など）、並びに/又は前述のいずれかの任意の組み合わせを含む。いくつかの実施形態において、曝露条件は、曝露期間、化合物の濃度、又は曝露期間及び化合物の濃度の組み合わせである。

20

## 【0489】

いくつかの実施形態において、共変量は、1つ以上の細胞（例えば、ペルターバゲン）において細胞状態遷移及び/又は摂動シグネチャを誘導する1つ以上の細胞に適用される化合物である。

## 【0490】

いくつかの実施形態において、共変量は、複数の細胞構成要素における細胞構成要素、又は第2の複数の細胞における細胞と関連付けられた知識用語（例えば、注釈）である。例えば、いくつかの実施形態において、共変量は、ゲノムワイド関連研究（GWAS）注釈、遺伝子セット濃縮アッセイ（GSEA）注釈、遺伝子オントロジー注釈、機能的及び/若しくはシグナル伝達経路注釈、並びに/又は細胞シグネチャ注釈である。いくつかの実施形態において、共変量は、NIH遺伝子発現オムニバス（GEO）、EBI Array Express、NCBI、BLAST、EMBL-EBI、GenBank、Ensembl、KEGG経路データベース、及び/又は任意の疾患特異的データベースを含むが、これらに限定されない、当該技術分野で既知の任意の公知の知識データベースから得られる。いくつかの実施形態において、共変量は、摂動（例えば、小分子）誘導遺伝子発現シグネチャを提供するデータベース、例えば、Library of Integrated Network-based Cellular Signatures (LINCS) L1000データセットから得られる。例えば、Duan, 2016, "L1000 CDS<sup>2</sup>: An ultra-fast LINCS L1000 Characteristic Direction Signature Search Engine," Systems Biology and Applications 2, article 16015を参照されたく、これは参照によりその全体が本明細書に組み込まれる。

30

40

## 【0491】

いくつかの実施形態において、複数の共変量は、少なくとも3個、少なくとも5個、少なくとも10個、少なくとも15個、少なくとも20個、少なくとも30個、少なくとも

50

40個、少なくとも50個、少なくとも60個、少なくとも70個、少なくとも80個、少なくとも90個、少なくとも100個、少なくとも200個、少なくとも300個、少なくとも400個、少なくとも500個、少なくとも600個、少なくとも700個、少なくとも800個、少なくとも900個、少なくとも1000個、少なくとも2000個、又は少なくとも3000個の共変量を含む。いくつかの実施形態において、複数の共変量は、5000個以下、1000個以下、500個以下、200個以下、100個以下、50個以下、又は20個以下の共変量を含む。いくつかの実施形態において、複数の共変量は、3~10個、10~50個、20~500個、200~1000個、又は1000~5000個の共変量を含む。いくつかの実施形態において、複数の共変量は、3個以上の共変量から始まり、5000個以下の共変量で終わる別の範囲内にある。

10

## 【0492】

いくつかの実施形態において、複数の共変量における共変量の各々は、細胞状態遷移及び/又は摂動シグネチャを誘導する1つ以上の細胞に適用される化合物であり、複数の共変量は複数の化合物である。いくつかの実施形態において、複数の共変量は、上記「化合物」と題されるセクションに開示されるように、複数の化合物からなる。

## 【0493】

図5は、細胞構成要素カウントデータ構造（例えば、目的の生理学的状態を通知する複数の共変量を集合的に表す第2の複数の細胞を使用して得られる）と、複数の細胞構成要素又はその表現を共通次元として使用する潜在表現とを組み合わせることによって形成される例示的な活性化データ構造を示す。これを達成するために、いくつかの実施形態において、第2の複数の細胞についてのカウントマトリックス（例えば、図4に示される第1の複数の細胞についてのカウントマトリックスと構造が類似している）及び潜在表現が一緒に乗算され、潜在表現マトリックスの重みがカウントマトリックスの正規化されたカウントによって乗算されるようにする。一般に、2つのマトリックスは、共通次元（例えば、第1のマトリックスのx軸及び第2のマトリックスのy軸）によって一緒に乗算される。第1及び第2のマトリックスのそれらの共通次元によるマトリックス乗算は、第1のマトリックス及び/若しくは第2のマトリックスに代替的に、又はそれに加えて、訓練されていないか、又は部分的に訓練されたモデルに適用することができる補助データの第3のマトリックスをもたらす。

20

## 【0494】

したがって、いくつかのそのような実施形態において、カウントマトリックスは、次元  $n_{\text{細胞}} \times n_{\text{遺伝子}}$  を有し、潜在表現は、次元  $n_{\text{遺伝子}} \times n_{\text{モジュール}}$  を有し、 $n_{\text{細胞}}$  は、第2の複数の細胞における細胞の数であり、 $n_{\text{遺伝子}}$  は、複数の細胞構成要素における細胞構成要素（例えば、遺伝子）の数、又はその表現であり、 $n_{\text{モジュール}}$  は、複数の細胞構成要素モジュールにおけるモジュールの数である。これは、カウントマトリックスにおける細胞構成要素の存在量を、各細胞（例えば、目的の1つ以上の共変量に対応する）がそのモジュール活性化によって特徴付けられ、得られたマトリックス表現（例えば、活性化データ構造）が、（例えば、 $n_{\text{遺伝子}}$  の共通次元を乗算した後の）次元  $n_{\text{細胞}} \times n_{\text{モジュール}}$  を有する空間にマッピングする。

30

## 【0495】

例えば、マトリックス乗算を使用する潜在表現及び細胞構成要素カウントデータ構造の組み合わせ、並びにマトリックス形態での結果として生じる活性化データ構造は、図5にまとめて示される。潜在表現（図5の上部パネルに示される）は、次元  $Z \times K$  を有し、ここで、 $Z$  は、細胞構成要素の数又はその表現であり、 $K$  は、細胞構成要素モジュールの数である。細胞構成要素カウントデータ構造（左下のパネルに示される）は、次元  $G \times Z$  を有し、ここで、 $G$  は、第2の複数の細胞における細胞の数であり、潜在表現に関して、 $Z$  は、細胞構成要素の数又はその表現である。 $Z$ （細胞構成要素の数又はその表現）を共通次元として使用するマトリックス乗算による組み合わせは、次元  $G \times K$  を有する結果として生じる活性化データ構造を生成する。それぞれの行の各々におけるそれぞれの列の各々についての各エントリは、それぞれの列に対応する第2の複数の細胞におけるそれぞれ

40

50

の細胞におけるそれぞれの細胞構成要素モジュールの各々の活性化を示す活性化重みである。したがって、図5に示されるように、モジュール1に対応するカウントは、細胞1に対応する活性化重み $\mu_{1-1}$ 、細胞Gに対応する活性化重み $\mu_{1-G}$ などを含む。

【0496】

いくつかの実施形態において、活性化データ構造における複数の活性化重みは、差次的モジュール活性化を含む。いくつかの実施形態において、差次的モジュール活性化（例えば、活性化データ構造における第2の複数の細胞における細胞間のそれぞれのモジュールの差次的活性化重み）は、関数 $(\mu_{i-1} - \mu_{i-2}) / (\text{var}_{i-1} + \text{var}_{i-2})^{-0.5}$ を使用してV-スコアを計算することによって得られ、 $\mu_{i-1}$ は、それぞれの条件i（例えば、共変量i）を有する細胞にわたるモジュール活性化の手段を示し、 $\text{var}_{i-1}$ は、条件iにおけるモジュール活性化の分散を示す。V-スコアは、分母内の細胞の数によって正規化されないt-スコアとして説明することができる。

10

【0497】

いくつかの実施形態において、活性化データ構造における第2の複数の細胞におけるそれぞれの細胞の各々は、それぞれの共変量を表す。いくつかの実施形態において、活性化データ構造における第2の複数の細胞におけるそれぞれの細胞の各々は、細胞状態遷移及び/又は摂動シグネチャを誘導する1つ以上の細胞に適用されるそれぞれの化合物を表す。

【0498】

したがって、いくつかの実施形態において、活性化データ構造は、第2の複数の細胞によって表される複数の化合物における各化合物への曝露に対応する（例えば、相関する、及び/又はそれに応答する）それぞれの細胞構成要素モジュールの活性化（例えば、活性化のレベル又は程度）を示す。例えば、第2の複数の細胞におけるそれぞれの細胞の各々がそれぞれのペルターバゲン（例えば、1つ以上の細胞が曝露される化合物並びに/又は細胞状態遷移及び/若しくは摂動シグネチャを誘導する化合物）を表すいくつかの実施形態において、活性化データ構造は、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの化合物による処置に相関する及び/又はそれに応答して、それぞれの細胞構成要素モジュールの活性化（例えば、誘導及び/又は差次的発現）を示すそれぞれの活性化重みを含む。

20

【0499】

いくつかの実施形態において、候補細胞構成要素モデルは、上記「モデルアーキテクチャ」と題されるセクションに記載されるように、本明細書に開示されるモデルアーキテクチャのいずれをも含む。

30

【0500】

いくつかの実施形態において、候補細胞構成要素モデルは、オートエンコーダ、スパースオートエンコーダ、及び/又はスパースマルチ読み出し、知識結合オートエンコーダである。いくつかの実施形態において、候補細胞構成要素モデルは、半教師ありモデルである。いくつかの実施形態において、候補細胞構成要素モデルは、1層ニューラルネットワーク（例えば、SoftMax及び/又はロジスティック回帰モデル）である。いくつかの実施形態において、候補細胞構成要素モデルは、次元Huber Outlier Regressorモデルである。

40

【0501】

いくつかの実施形態において、候補細胞構成要素モデルは、複数の層を含むスパースマルチ読み出し、知識結合オートエンコーダであり、第1の層は潜在表現を得るために使用され、第2の層は細胞構成要素モジュール知識構築物（例えば、共変量重みマトリックス）を得るために使用される。

【0502】

いくつかの実施形態において、候補細胞構成要素モデルを訓練することは、マルチタスク策定におけるカテゴリ交差エントロピー損失を使用して実施され、複数の共変量における共変量の各々が、複数のコスト関数におけるコスト関数に対応し、複数のコスト関数に

50

おけるそれぞれのコスト関数の各々が、共通の重み付け係数を有する。

【0503】

いくつかの実施形態において、候補細胞構成要素モデルを訓練して、目的の生理学的状態と関連付けられた細胞構成要素モジュールのセットにおける第1の細胞構成要素モジュールを識別するようにモデルを訓練する。モデルを訓練するための方法については、本明細書に更に詳細に記載される。本明細書に開示される方法及び/又は実施形態のいずれかは、上記「モデル訓練」と題されるセクションに記載されるように、候補細胞構成要素モデルの訓練において使用することが企図される。

【実施例】

【0504】

V. 実施例

本明細書に提供されるのは、化合物を生理学的状態と関連付けるためのモデルの例示的な性能尺度及び治療上の適用である。

【0505】

実施例1. 脂肪酸関連細胞プロセスの活性化のための化学構造の予測。

この実施例では、細胞構成要素モジュールを最初に定義した。これは、細胞が目的の生理学的状態と関連付けられた異なる状態を表す細胞についての発現データを得ることによって行われた。これは、当初出願された請求項27を追跡する。細胞構成要素存在量値は、細胞の各々から測定され、このデータは、細胞構成要素をクラスタ化するために使用される。細胞によって表される様々な状態にわたって発現値が互いに相関しているこれらの細胞構成要素を、細胞構成要素モジュールにグループ化する。これは、いくつかの細胞構成要素モジュールをもたらし、それらの各々は、細胞構成要素試料の異なるサブセットを含む。いくつかの実施形態において、細胞構成要素モジュールの各々は、細胞構成要素の異なるサブセットを有するが、ある細胞構成要素モジュールにおける細胞構成要素と別の細胞構成要素モジュールにおける細胞構成要素との間に重複がある可能性がある。

【0506】

更に、この例では、追加の訓練データが、第2の訓練セットの形態で得られる。この第2の訓練セットはまた、細胞構成要素についての単一細胞存在量データを含む。しかしながら、この第2の訓練セットにおいて、各細胞は、複数の訓練化学化合物において異なる化学化合物に曝露されている。この訓練セットにおいて、既知の量は、それぞれ異なる化学化合物のフィンガープリントであり、そのような化合物に曝露された細胞の得られる細胞構成要素存在量データである。第2のデータセットについてのデータは、細胞構成要素同一性のための第1の軸、及び細胞同一性のための第2の軸を有する、カウントマトリックス502(図5に示される)として配置することができる。したがって、カウントマトリックス502における各エレメントは、所与の細胞内の所与の細胞構成要素の存在量である。更に、(特定の細胞に対応する)カウントマトリックス502におけるそれぞれの列の各々は、特定の細胞が曝露された特定の化合物で標識される。したがって、カウントマトリックス502の各列は、特定の化合物(例えば、訓練化合物)で標識されるが、各エレメントは、対応する細胞(X軸)についての対応する細胞構成要素(Y軸)のカウントである。

【0507】

図5に示されるように、第1のデータセット(潜在表現404)及び第2のデータセット(カウントマトリックス502)からのデータは、組み合わせられて、活性化データ構造(例えば、図5に示されるような活性化データ構造504)を形成する。例えば、これを達成する1つの方法は、第1の軸が細胞構成要素モジュールを表し、第2の軸が細胞構成要素の各々を表すように、潜在表現404における行として細胞構成要素モジュールを配置することである。このようにして、活性化データ構造504を生成するために、潜在表現404及びカウントマトリックス502は、マトリックス乗算を介して、それらの共通軸、細胞構成要素の数によって乗算されて、活性化データ構造504に到達する。活性化データ構造504は、カウントマトリックス502からの細胞同一性軸及び潜在表現50

10

20

30

40

50

4からの細胞構成要素モジュール軸を保持する。異なる細胞型に対して異なる活性化構造を形成することができる。すなわち、カウントマトリックス502を形成するために使用される細胞は、目的の特定の疾患状態を表すことができる。したがって、異なる疾患状態又は目的の他の表現型について、異なる活性化データ構造504を形成することができる。

#### 【0508】

図6を参照すると、いくつかの例では、活性化データ構造504の各行(図5から、現在は図6の上部にある)は、異なるモデル601についての訓練データとして機能する。例えば、モデル601が行604-1の重み(重み $1-1$ から重み $1-w$ )を含み、化合物1からWがそれぞれ細胞構成要素モジュール1を活性化する程度を表す場合を考慮する。このモデル601は、活性化データ構造504の行640のエレメントについて訓練され、これは、訓練化合物1、...、Gの各々が細胞構成要素モジュール1を活性化する程度を提供する。この訓練では、まず、細胞1が曝露された化合物のフィンガープリント表現がモデル601に入力される。この入力に回答して、細胞構成要素モジュール1についてのモデル601は、Predと呼ばれる活性化値を出力する。図6の命名法における値 $1$ 。この出力活性化値は、活性化データ構造504のAct $1-1$ である、実際の活性化値と比較される。次に、細胞2が曝露された化合物のフィンガープリント表現がモデル601に入力される。この入力に回答して、モデルは活性化値(Pred、値 $2$ )を出力する。この出力活性化値は、活性化データ構造504のAct $1-2$ である、化合物2についての実際の活性化値と比較される。このプロセスは、細胞Gを通じて進行する。細胞Gが曝露された化合物のフィンガープリント表現は、モデル601に入力される。これに回答して、モデルは活性化値(Pred、値 $G$ )を出力する。この出力活性化値は、活性化データ構造504のAct $1-G$ である、細胞Gについての実際の活性化値と比較される。この例では、W及びGは同じ値を有する。このようにして、細胞構成要素モジュール1について図5に概説されるように、活性化データ構造を導出するために使用される化合物の訓練セットにおける各化合物について結果として生じる予測(Pred、値)が存在する。(活性化値の)上述の計算された予測を、これらの化合物の各々について上述の実際の活性化値と比較し、予測された活性化値と実際の活性化値との間の差を使用して、逆伝搬及び関連するモデル改良技術を使用してモデル601を更に訓練する。

#### 【0509】

したがって、結果は、細胞構成要素モジュールの各々のうちの1つである、一連の訓練されたモデル601である。試験化合物のフィンガープリントは、訓練されたモデルの各々に入力され得、それぞれの訓練されたモデル601の各々は、予測された活性化値を出力し、その大きさは、それぞれの訓練されたモデルに対応する細胞構成要素モジュールが試験化合物によって活性化されるかどうかを示す。ここで、プロセスの概要が説明されたので、ステップの各々は、この例で使用される実験データと併せて説明される。

#### 【0510】

以下のプロセスによって、第1の細胞構成要素モジュール(図1、図4 132-1)を識別する。電子形式で1つ以上の第1のデータセットを得る。1つ以上の第1のデータセットは、複数の注釈付き(例えば、標識されたか、又は既知の)細胞状態を集合的に表す第1の複数の細胞(例えば、20個以上の細胞)に関するデータを含む。第1のデータセットは、第1の複数の細胞におけるそれぞれの細胞の各々について、複数の細胞構成要素(例えば、10個以上の細胞構成要素)におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含む。例えば、各細胞についての転写データ。このようにして、複数のベクトルがアクセスされるか、又は形成される。複数のベクトルにおけるそれぞれのベクトルの各々は、複数の構成要素におけるそれぞれの細胞構成要素に対応し、対応する複数のエレメントを含む。ベクトルの対応する複数のエレメントにおけるそれぞれのエレメントの各々は、第1の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を表す対応するカウントを有する。したがって、いくつかのそのような実施形態において、複数の細胞状

態における細胞状態の各々についての転写データが得られる。

【0511】

例示するために、図4に例示される形態のカウントマトリックス402が形成される。この例では、前脂肪細胞において代謝活性プロセスを誘導することが知られている小分子ペルターバゲンを使用した。前脂肪細胞株のアリコートにペルターバゲンに24時間曝露し、摂動状態における細胞株の曝露アリコートについてscRNA-seq読み出しを得た。scRNA-seq読み出しは、ペルターバゲンに曝露されていない細胞株のアリコートについても得られ、これらの読み出しは、対照条件を表している。このようにして、図14Aのブロック1504に従って、第1のデータセットを得、これは、第1の複数の細胞におけるそれぞれの細胞の各々について、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含み、それによって、複数のベクトルにアクセスしたか、又はそれを形成した。すなわち、ペルターバゲンに曝露された細胞及びペルターバゲンに曝露されなかった細胞（対照細胞）の両方で測定された各細胞構成要素（例えば、遺伝子）の発現値は、図4に例示されるカウントマトリックス402の要素を形成した。図4に例示され、図14Aのブロック1510に記載されるように、カウントマトリックス402は、各細胞構成要素についてのベクトルを含み、したがって、複数のベクトルが存在する。複数のベクトルにおけるそれぞれのベクトルの各々は、(i)複数の構成要素におけるそれぞれの細胞構成要素に対応し、(ii)対応する複数の要素を含む。

10

【0512】

例えば、細胞構成要素1（例えば、遺伝子1）について、カウント1-1、...、カウント1-Nは、細胞1からNにおける遺伝子1の発現の測定値であり、N細胞のいくつかは、ペルターバゲンに曝露されており、いくつかは、ペルターバゲンに曝露されておらず、これらのカウントは、細胞構成要素1についてのベクトルの要素を形成する。すなわち、図14Aのブロック1512に従って、細胞構成要素1のベクトルの対応する複数の要素におけるそれぞれの要素の各々は、第1の複数の細胞におけるそれぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を表す対応するカウントを有する。この例は、2つの状態（ペルターバゲンに曝露されたか、又はされていない）を含むが、原則として、異なる濃度のペルターバゲン、曝露時間などの、任意の数の状態を包含することができる。

20

30

【0513】

図14Aのブロック1514によれば、この実施例1には、対照（ペルターバゲンへの曝露なし）及びペルターバゲンの曝露という2つの注釈付きの状態がある。すなわち、複数の注釈付きの細胞状態における1つの注釈付きの細胞状態は、曝露条件（例えば、曝露期間、ここでは24時間）下での化合物（ここでは、ペルターバゲン）への第1の複数の細胞における細胞の曝露である。この例は、2つの状態（ペルターバゲンに曝露されたか、又はされていない）からなるが、原則として、異なる濃度のペルターバゲン、曝露時間などの、任意の数の状態を包含することができる。

【0514】

カウントマトリックス402は、フィルタリング及び正規化ステップを介して前処理され、高いシグナル対ノイズ比を有するいくつかの遺伝子を含む前処理されたカウントマトリックスをもたらした。

40

【0515】

複数のベクトルを使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別する。複数の細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々は、複数の細胞構成要素のサブセットを含む。複数の細胞構成要素モジュールは、(i)複数の候補細胞構成要素モジュール、及び(ii)複数の細胞構成要素又はその表現によって次元決定された潜在表現で配置され、複数の細胞構成要素モジュールは、10を超える細胞構成要素モジュールを含む。

【0516】

50



いくつかの実施形態において、候補細胞構成要素モジュールの各々は、候補転写フィンガープリントである。

【0517】

この実施例では、カウントマトリックス402を使用して、細胞構成要素モジュール132を識別した。これは、図14Bのブロック1526に従って行われた。複数のベクトル（図4のカウントマトリックス402の各行）を使用して、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々を識別し、複数の候補細胞構成要素モジュールにおける候補細胞構成要素モジュールの各々は、複数の細胞構成要素のサブセットを含む。

【0518】

これにより、(i)複数の候補細胞構成要素モジュール、及び(ii)複数の細胞構成要素、又はその表現によって次元決定された潜在表現が得られ、複数の候補細胞構成要素モジュールは、10を超える細胞構成要素モジュールを含む。この潜在表現の例は、図4の潜在表現404であり、それぞれの候補細胞構成要素モジュール132の各々について、どの細胞構成要素がそれぞれの候補細胞構成要素モジュール内にあるかを示す。

【0519】

潜在表現404は、図14Cのブロック1528に従って形成され、複数のベクトル（カウントマトリックス402の細胞構成要素ベクトル）を使用して、複数のベクトルにおけるベクトルの各々の対応する複数のエレメントの各々を使用して複数のベクトルに相関モデルを適用することによって、複数の候補細胞構成要素モジュール（潜在表現404の）における候補細胞構成要素モジュールの各々を識別した。特に、相関ベースのコスト関数が最適化され、これは、細胞構成要素ベクトル間の最近傍関係を定義する最近傍グラフを計算し、カウントマトリックス402の細胞構成要素ベクトル間の相関を計算することに相当した。複数の細胞にわたって互いに高い相関を有する細胞構成要素（ここでは遺伝子）は、最終的に最近傍となり、ライデン（Leiden）又は任意の他のグラフクラスタリング方法を使用してグラフをクラスタリングすることによって、潜在表現402内に細胞構成要素モジュールを形成した。スパースオートエンコーダコスト関数を最適化することは、pytorch又はtensorflowに実装されているような標準訓練を使用して、その重みのL0正規化、及び再構築損失を伴う1層オートエンコーダの訓練に相当した）。この実施例では、これは、訓練中に108個の細胞構成要素モジュールが学習されることをもたらした。すなわち、図4の潜在表現404は、108個の細胞構成要素モジュール132を有し、各々が、カウントマトリックス402において発現データが利用可能であった細胞構成要素の独立したサブセットを有する。

【0520】

108個の細胞モジュールのうち、「モジュール78」と称される細胞構成要素モジュール132は、摂動試料及び対照試料にわたって計算された細胞構成要素の各々についてのtスコアを平均化するとき、最も強い活性化を示した。言い換えれば、カウントマトリックスデータ内の発現データを使用して、潜在表現404内のそれぞれの細胞構成要素モジュールの各々について、ペルターバゲンに曝露された細胞とペルターバゲンに曝露されていない細胞との間のそれぞれの細胞構成要素モジュールにおける細胞構成要素の各々の差次的発現に関するtスコアを実施することによって、細胞構成要素を検証した。更に、モジュール78は、脂肪酸及び脂質に関連する生物学的プロセスに関与する細胞構成要素で濃縮されている。要約すると、モジュール78は、代謝活性のマーカーである、FABP3を含む、28個の遺伝子からなる。

【0521】

細胞構成要素モジュールに加えて、細胞を訓練化合物に曝露したときの細胞ベースの細胞構成要素応答データが必要である。

【0522】

したがって、1つ以上の第2のデータセットを電子形式で得た。1つ以上の第2のデータセットは、第2の複数の細胞からのデータを含む。第2の複数の細胞は、20個以上の

10

20

30

40

50

細胞を含む。第2の複数の細胞は、目的の生理学的状態を通知する複数の共変量を集合的に表した。例えば、複数の共変量は、いくつかの例では、訓練化合物である。次に、第2の複数の細胞における細胞の各々について、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量が取得され、それによって、(i)第2の複数の細胞、及び(ii)複数の細胞構成要素又はその表現によって次元決定された細胞構成要素カウントデータ構造が得られる。

【0523】

これは、図14Cのブロック1538に従っており、第2のデータセットが電子形式で得られた状態は、第2の複数の細胞におけるそれぞれの細胞の各々について、第2の複数の細胞が、20個以上の細胞を含み、目的の生理学的状態を通知する複数の共変量(ここでは複数の異なる化学化合物)を集合的に表し、複数の細胞構成要素におけるそれぞれの細胞構成要素の各々について、それぞれの細胞におけるそれぞれの細胞構成要素の対応する存在量を含み、それによって、(i)第2の複数の細胞、及び(ii)複数の細胞構成要素又はその表現によって次元決定された細胞構成要素カウントデータ構造を得る。

10

【0524】

このカウントマトリックスの形態の例示は、図5のカウントマトリックス502である。図5のカウントマトリックス502に示されるように、それぞれの細胞構成要素(例えば、遺伝子)の各々について、第2の複数の細胞における細胞の各々についての発現データが存在する。例えば、複数の遺伝子の各々の転写活性を、第2の複数の細胞にわたって測定する。細胞の各々は、共変量、ここでは訓練化学化合物に曝露されている。

20

【0525】

複数の細胞構成要素又はその表現を共通次元として使用して、細胞構成要素カウントデータ構造及び潜在表現を組み合わせることによって活性化データ構造を形成し、活性化データ構造は、複数の細胞構成要素モジュールにおける細胞構成要素モジュールの各々について、第2の複数の細胞における細胞の各々について、それぞれの活性化重みを含む。

【0526】

カウントマトリックス502は、図5に示される活性化データ構造504を得るために、潜在表現404によって乗算されたマトリックスであった。活性化データ構造504は、それぞれの細胞構成要素モジュールの各々について、第2の複数の細胞における細胞の各々について、活性化値  $Act_{k-g}$  を有し、その値は、カウントマトリックス502による潜在表現404の対応するマトリックス乗算によって決定される。

30

【0527】

(i)活性化データ構造を候補モデルに入力したときに、活性化データ構造内に表される細胞構成要素モジュールの各々における複数の共変量における各共変量の不在又は存在の予測と、(ii)細胞構成要素モジュールの各々における各共変量の実際の不在又は存在との間の差を使用して、候補細胞構成要素モデルを訓練し、訓練することは、差に応答して、候補細胞構成要素モデルと関連付けられた複数の共変量重みを調整する。

【0528】

活性化データ構造502は、それ自体が次元N化合物×M細胞構成要素モジュールの潜在表現602である、図6のモデル601についての訓練データ(標識データ)として機能した。この実施例では、8000個の異なる化合物及び108個の細胞構成要素モジュールを考慮した。したがって、図5の命名法では、Zは108であり、Gは8000であった。活性化データ構造は、2つの方法で訓練及び試験セットに分割された。まず、1200個の化合物を試験セットにグループ化し、残りの6800個の化合物を訓練セットにグループ化した「ランダム分割」を選択した。また、「クロス骨格分割」は、試験セットが訓練セットとは異なる骨格を有する化合物を含むことを保証するオープンソースソフトウェアパッケージRDKitの機能を使用して定義された。

40

【0529】

図6に示されるように、活性化データ構造504のそれぞれの行の各々は、それぞれの行によって表される対応する細胞構成要素モジュールの細胞構成要素を誘導する可能性が

50

高い化合物を表すベクトルである。モデル 601 の各インスタンスは、活性化データ構造 504 の行で訓練された。活性化データ構造 504 は、6800 個の訓練化合物を使用して形成された。所与のモデル 601 について、特定の化学化合物のフィンガープリントがモデル 601 に入力され、この入力に応答して、対応する細胞構成要素モジュールについての予測された活性化値が計算される。この予測された活性化値は、活性化データ構造 504 内の対応するエレメントにおける実際の活性化値と直接比較され得る。したがって、このようにして、(i) モデル 601 への活性化データ構造 504 の入力時に活性化データ構造 504 に表される細胞構成要素モジュールの各々についての訓練化合物における各化合物の不在又は存在の予測と、(ii) 細胞構成要素モジュールの各々についての各化合物の実際の不在又は存在との間の差を計算し、差に応答して候補細胞構成要素モデルと関連付けられた複数の共変量重み 604 を調整することによってモデル 601 を訓練するために使用することができる。図 6 に示されるように、複数の共変量重みは、複数の細胞構成要素モジュールにおけるそれぞれの細胞構成要素モジュールの各々について、それぞれの共変量の各々について、それぞれの共変量が、活性化データ構造にわたって、それぞれの細胞構成要素モジュールと相関するかどうかを示す対応する重みを含む。いくつかの実施形態において、細胞構成要素モジュールの各々について異なるモデル 601 が存在した。言い換えれば、図 6 を参照すると、いくつかの実施形態において、各行 604 は異なるモデル 601 内にある。したがって、そのような実施形態において、そのようなモデル 601 の各々は、活性化データ構造における対応する行（例えば、それぞれのモデル 601 と同じ細胞構成要素モジュールに対応する行）を使用して訓練される。

10

20

#### 【0530】

図 6 に示すように、訓練されたモデル 601（又は複数のモデル）は、各共変量（ここでは、訓練化学組成物）についての重みを提供する。すなわち、モデル 601 の潜在表現 602 は、各共変量（化学組成物）が細胞構成要素モジュールの活性化にどの程度関連するかを説明する重み（例えば、図 6 の重み  $w_{1-1}$  又は行 604 - 1）を提供する。そのような重みは、細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれの数値的活性化スコアとみなされる。細胞構成要素モジュールの各々についての異なるモデル 601 が形成される実施形態において、潜在表現 602 は、各モデル 601 の集合潜在表現である。いくつかの実施形態において、表現の各重みは分類的である（例えば、化合物は細胞構成要素モジュール「0」に影響を及ぼすか、又は化合物は細胞構成要素モジュール「1」に影響を及ぼさない。他の実施形態において、各重みは、連続スケール上にあり、スケールの一端は、訓練化合物が細胞構成要素モジュールに大きく影響することを示し、スケールの他端は、訓練化合物が細胞構成要素モジュールに影響しないことを示す。本明細書で使用される場合、「影響する」という用語は、用途依存性であるが、概して、化合物の不在又は存在が、細胞構成要素モジュールにおける細胞構成要素の存在量を変化させることを意味する。

30

#### 【0531】

モデル 601 の訓練のために、この例では、図 6 の活性化データ構造 504 に表される化合物の SMILES 表現は、ECFP4 フィンガープリント表現、更にはグラフ表現に変換される。その後、2つのモデルが訓練される。すなわち、モデル 601 は、この例では、2つの異なるモデルのアンサンブルである。A) 完全に接続されたニューラルネットワークアーキテクチャを使用して ECFP4 表現を訓練し、B) メッセージパッシングニューラルネットワーク (MPNN) を使用してグラフ表現を訓練する。この訓練を実施するために、オープンソースソフトウェアパッケージ `pytorch` 及び `DGL` を使用した。訓練されていないモデル 601 は、訓練セットにおけるそれぞれの化合物の各々のそれぞれの化学構造の各々について、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々について、(i) それぞれの化合物の化学構造のフィンガープリントを訓練されていないモデルに入力したときのそれぞれの細胞構成要素モジュールについてのそれぞれの計算された活性化スコアと、(ii) 細胞構成要素モジュールのセットにおけるそれぞれの化合物についてのそれぞれの細胞構成要素モジュールのそれぞれ

40

50

の数値的活性化スコア（活性化データ構造504から得た）との間のそれぞれの差を使用して訓練され、訓練することが、差に応答して訓練されていないモデル601と関連付けられた複数のパラメータを調整し、複数のパラメータが、100以上のパラメータを含み、それによって、訓練されたモデルを得る。

#### 【0532】

上述のように、この例では、モデル601は、(i)SMILES文字列の標準フィンガープリント上の完全に接続されたネットワークであって、ネットワークアーキテクチャは、ReLU活性化を有する3層ネットワークである、完全に接続されたネットワーク、及び(ii)DGLライブラリからのMPNNネットワークのアンサンブルである。化学構造情報を入力すると、モデル601は、それが訓練された細胞構成要素モジュール132の各々の活性化スコアを提供する。

10

#### 【0533】

実際に、いくつかの実施形態において、この例では、細胞構成要素モジュールの各々のための別個のアンサンブルモデル601が存在する。言い換えれば、モデル601は、化学構造の入力時に複数の細胞構成要素モジュールの各々に対して別個の活性化スコアを提供するマルチタスクエンコーダであった。なお更に、いくつかの実施形態において、上で説明したように、それぞれの細胞構成要素モジュールの各々について別個のモデル601が存在する。そのような実施形態において、そのようなそれぞれのモデル601の各々は、対応する細胞構成要素モデルに対する化合物の各々についての活性化重みを含む。

#### 【0534】

現在訓練されているそれぞれのモデル601の各々は、訓練セットの一部であるか否かにかかわらず、任意の化合物について、その対応する細胞構成要素モジュールについての活性化スコアを提供する。すなわち、各モデル601は、その対応する細胞構成要素モジュールが試験化合物と関連付けられているかどうかを報告することができる。それがあある場合、モデルは、その対応する細胞構成要素モジュールが試験化合物と関連付けられていることを示すスコアを出力する。いくつかの実施形態において、このスコアは分類的である（例えば、対応する細胞構成要素モジュールが試験化合物と関連付けられている場合は「1」であり、関連付けられていない場合は「0」である）。いくつかの実施形態において、このスコアは、例えば、1に近い数（例えば、0.85）が、対応する細胞構成要素モジュールが試験化合物と関連付けられている可能性を示す、0~1のスケールでの確率又は尤度である。いくつかの実施形態において、このスコアは、「A」から「B」の連続スケール上にあり、A及びBは2つの異なる数である。各々が異なる細胞構成要素モジュールに対応するいくつかのモデル601が存在するので、試験化合物をいくつかの異なるモデル601に対して実行して、どの細胞構成要素モジュールが化合物によって活性化されるか（それと関連付けられるか）を決定する。各例において、化学構造は、上で説明したようにフィンガープリントに変換され、各モデルに適用されるのはこのフィンガープリントである。生物学的観点から、所与の試験化合物は、任意の数の異なる細胞構成要素モジュール（例えば、1、2、3、4、5、又はそれよりも多い）を活性化し得ることが予想され得ることに留意されたい。更に、本開示に記載のアプローチは、モデル601が訓練されていないが、どの細胞構成要素モジュールが試験化合物によって活性化されるべきかが知られている化合物を試験することによって検証することができる。これは、以下に示すようにこの例で行われた。特に、化合物を生理学的状態と関連付けるための訓練されたモデル601は、この例では4倍に検証された。この試験は、当初出願された請求項1を追跡する。

20

30

40

#### 【0535】

まず、モデル601からのモデル予測は、ハイスループットスクリーニングから上記の1200個のランダムに選択された目に見えない化合物によって、また6800個の化合物訓練セットに対して上記の1200個の重複しない骨格を有する化合物によって誘導された脂肪酸生成関連細胞構成要素モジュールの活性化について得られた。ランダムに選択された化合物について得られたそれぞれのモデル601の予測（予測された細胞構成要素

50

活性化スコア)を図10Bに示す。すなわち、図10Bは、2つの異なるモデル601、すなわち、1つは細胞構成要素モジュール78「モジュール78」、1つは細胞構成要素モジュール「90」からの結果を示す。モジュール78は、細胞代謝にとって重要な脂肪酸関連細胞プロセスを表し、その対応する訓練されたモデル601は、高い決定係数を示した( $R^2 = 0.28$ )。対照的に、同じscRNA-seqデータセットから学習した細胞構成要素「モジュール90」のための訓練されたモデル601は、細胞代謝とは無関係であり(モジュール90における細胞構成要素は脂肪酸関連プロセスとは関連しない)、低い決定係数を有した( $R^2 = 0.08$ )。全てのベンチマークは、非常に有意な相関をもたらした(それぞれ、ピアソン相関係数 $p_s =$ 約0.5及び約0.2)。

#### 【0536】

当初出願された請求項1に記載の言語において、この第1の検証アプローチは、試験化学化合物(ハイスループットスクリーニングから、また6800個の化合物訓練セットに対する上記の1200個の重複していない骨格を有する化合物によって、記載された1200個のランダムに選択された目に見えない化合物のうちの一つ)を目的の生理学的状態(ここでは、この例では、細胞代謝にとって重要な脂肪酸関連細胞プロセス)と関連付ける方法を提供する。この方法は、メモリ及び1つ以上のプロセッサを含むコンピュータシステムにおいて、試験化学化合物の化学構造のフィンガープリントを得ることを含む。したがって、試験化学化合物の化学構造のフィンガープリントが得られ、それはこの例で図1の各モデル601に入力されるものである。当初出願された請求項1の文脈において、モデルは、モデルと称される。このモデルは、アンサンブルモデルを包含し、アンサンブルモデルにおける各コンポーネントモデルは、図6のモデル601について列挙されたパラメータの単一の行を含み、行は、コンポーネントモデルと関連付けられた所与の細胞構成要素モジュールについての重みに関するパラメータである。図6において、そのような重みは単一の行として表されるが、それらがアンサンブルモデルのコンポーネントモデルにおいて行の形式であるという要件はなく、その任意の等価物は本開示の範囲内であることを理解されたい。更に、図6のモデル601は、それが訓練された各化合物についての単一の重みを含み、これは回帰に基づいてモデル601に好適であるが、いくつかの実施形態において、モデル601の重みの数と、モデルが訓練された化合物の数との間に明確な関係は存在しない。いくつかの実施形態において、モデル601は、100以上、1000以上、10,000以上、又は100,000以上のパラメータを含む。

#### 【0537】

当初出願された請求項1によれば、試験化合物のフィンガープリントは、モデルに入力される。当初出願された請求項1に記載されているように、モデルは、100以上のパラメータを含む。言い換えれば、モデル出力の計算は、試験化合物のフィンガープリントを入力すると、精神的に実施することができない。モデルは、フィンガープリントのモデルへの入力に回答して1つ以上の計算された活性化スコアを出力する。1つ以上の計算された活性化スコアにおけるそれぞれの計算された活性化スコアの各々は、細胞構成要素モジュールのセットにおける対応する細胞構成要素モジュールを表す。この例では、モデルはモデル601のアンサンブルであり、各々が異なる細胞構成要素モジュールを表し、したがって、アンサンブル内の各モデル601は、細胞構成要素モジュールのセットにおける単一の対応する細胞構成要素モジュールを表す1つ以上の計算された活性化スコアにおける計算された活性化スコアを出力する。この点で、及び上述したように、細胞構成要素モジュールのセットにおけるそれぞれの細胞構成要素モジュールの各々は、複数の細胞構成要素の独立したサブセットを含む。更に、細胞構成要素モジュールのセットにおける少なくとも第1の細胞構成要素モジュールは、目的の生理学的状態と関連付けられる。この実施例では、モジュール78は、目的の生理学的状態と関連付けられる。図10Bに示されるように、モジュール78を正しく活性化し、したがって、モジュール78の目的の生理学的状態(細胞代謝に重要な脂肪酸関連細胞プロセス)と関連する化合物は、(例えば、第1の閾値基準を満たす第1の細胞構成要素モジュールについてのそれぞれの計算された活性化スコアによって)識別される。

10

20

30

40

50

## 【0538】

特許請求されるアプローチの第2の検証として、次に、モジュール78及び90についてのそれぞれの訓練されたモデルを、訓練中に図6のモデル601に導入されていない別の試験セットである前脂肪細胞に曝露された特定の小分子である「合成ヒット」のscRNA-seq特性評価に適用した。図10Dは、合成ヒットによるモジュール78についての訓練されたモデル601によって示される活性化の高い相関及び忠実な予測を、合成ヒットによるモジュール90についての訓練されたモデル601によって示されるほとんどない活性化、又は全くない活性化と比較して示す。

## 【0539】

第三に、モジュール78のための訓練されたモデル601を使用して、公開データベース内の500万個の化合物からサンプリングされた200,000個の化合物のランダムなサブセットについて、細胞構成要素モジュール78(モジュール78)についての細胞構成要素活性化スコアを予測した。このことから、細胞構成要素モジュール78を高度に活性化すると予測される上位50個の化合物を選択し、LINCS L1000データセットからの化合物及び本明細書において既知のピペリジン含有化合物(「KPC」)と称される既知の化合物の化学構造に由来する合成ヒット類似体を含むデータベース内の化合物のセットと比較した。この比較の分布を図10Eに示す。分布の末端で、細胞構成要素モジュール78についての訓練されたモデル601について得られた予測は、LINCS及び合成ヒットにおける全ての化合物を有意に上回る化合物を識別した。このアプローチは、特定の所望の細胞プロセスに対して化学構造を最適化するための方法を強調する。

## 【0540】

第4に、上位50の予測で識別された化学構造を視覚的に検査し、既知の脂肪組織標的化ファーマコフォア)を表す明白な化学構造を含み、したがって、モジュール78と関連付けられた細胞構成要素モジュールを正当に活性化することが見出された。

## 【0541】

この第1の例はまた、当初出願された請求項58を追跡する。請求項1と請求項58との間の違いは、細胞構成要素モジュールに対する摂動シグネチャの1つである。摂動シグネチャは、摂動にさらされた細胞とそうではない細胞の発現を比較することによって得られる。したがって、前脂肪細胞において代謝活性プロセスを誘導することが知られている低分子ペルターバゲンを使用することができる。前脂肪細胞株をペルターバゲンに24時間曝露し、scRNA-seq読み出しを摂動状態及び対照状態について得ることができる。このことから、摂動シグネチャを得ることができる。あるいは、別個の摂動シグネチャは、第2のデータセットに使用される化学共変量のうちのいずれか1つに曝露された細胞の細胞発現を比較することによって得ることができる。実際に、第2のデータセットに使用される化学共変量の各々について、このようにして別個の摂動シグネチャを得ることができる。そのような摂動シグネチャの各々は、そのような重みの各々が現在2値スケールではなく連続スケール上にあることを除いて、潜在表現404内の行の形態を有する。例えば、いくつかの実施形態において、各重みは、0から1(又はいくつかの他の範囲の「A」から「B」であり、A及びBは、-100及び100などの2つの異なる数字である)の間の連続スケール上の値である。そこから、訓練のプロセスは、潜在表現404、カウントマトリックス502、活性化データ構造、及びコンポーネントモデル601の訓練の使用に関して上述したものと同一であり、そのようなモデルの各々は、現在、摂動シグネチャのセットにおける異なる摂動シグネチャを表す。

## 【0542】

実施例2．胎児赤血球生成プログラムを活性化し、T細胞枯渇をブロックするための化学構造の予測。

2つの追加の例では、胎児赤血球生成及びT細胞枯渇に関連する2つのscRNA-seqデータセット上の2つのモデルを訓練した。

## 【0543】

胎児赤血球生成のために、CD34造血幹細胞をツール化合物CLT-AAA-12で

処理し、これに関して、胎児赤血球生成のエンドポイントマーカー、特にフローサイトメトリーによる読み出しとしてのアッセイにおけるF細胞の数が誘導されることが以前に確立されている。

【0544】

T細胞枯渇のために、ナイーブT細胞を枯渇誘導培地で処置した。

【0545】

両方の細胞系は、scRNA-seqで特徴付けられる。その後、薬物リフレクターモデル（参照により本明細書に組み込まれる、2019年7月15日に出願された「Methods of Analyzing Cells」と題された米国特許出願第16/511,691号を参照されたい）を、摂動細胞対照細胞によって定義される細胞状態遷移をそれらのそれぞれの試料に入力することによって、scRNA-seqデータセットに適用した。薬物リフレクターは、薬物リフレクター潜在表現における8000個の化合物の各々について細胞状態活性化スコアを割り当てる。これにより、両方の遷移（胎児ヘモグロビン及びT細胞枯渇）について細胞状態活性化スコアを有する2つのベクトルが生じる。これらの2つのベクトルは、モデル601についての訓練データとして機能する。

【0546】

このモデルを使用して、造血幹細胞における胎児赤血球生成を活性化する化合物及びT細胞枯渇を予測した。造血幹細胞における胎児赤血球生成は、近年、鎌状赤血球疾患に対する画期的なCRISPR療法につながった細胞プロセスであり、一方、T細胞枯渇は、がんに対するチェックポイント阻害剤療法のより広範な成功を妨げる重要な機構である。

【0547】

予測は、公開データベース内の500万個の化合物からサンプリングされた2,000個の化合物のサブセットを使用して実施され、サブセットは、ランダムに又は骨格上で分割された。図11の上部パネルは、ランダムに分割された、及び骨格上にある2,000個の化合物の試験セット上でのこの実施例のモデルの性能を示し、造血幹細胞における胎児赤血球生成に関連するヒット化合物CLT-AAA-12の摂動シグネチャとともに、サンプリングした化合物の有意な $R^2$ 及び相関係数 $p_s$ を示している。図11の下部パネルは、ランダムに分割された、及び骨格上にある2,000個の化合物の試験セットの性能を示し、T細胞枯渇に関連する細胞遷移シグネチャとともに、サンプリングした化合物の有意な $R^2$ 及び相関係数 $p_s$ を示している。したがって、図11は、モデル601が、目的の摂動シグネチャ及び/又は細胞遷移シグネチャと同じ細胞挙動効果を誘導する新しい骨格を予測することができることを実証する。

【0548】

実施例3．疾患クリティカル細胞挙動に基づく特徴属性：新しい分子の設計のためのファーマコフォアの予測。

実施例1に記載されるように、本明細書に開示されるシステム及び方法に従って予測される化学構造を使用して、目的の生理学的状態（例えば、脂肪組織標的化）に潜在的に関連している、ファーマコフォアなどの分子的特徴を識別することができる。実施例1と同様に、これらのファーマコフォアは、既知の化学構造によって検証することができるか、又は更なる検証のために新規の構造を提示することができる。例えば、ファーマコフォアに基づくアルゴリズムの例示的な使用事例は、Base of Bioisosterically Exchangeable Replacements (BoBER) データベースを含む、以前に文献に記載された機能的意味を有するファーマコフォアのデータベースを活用することを含む。使用事例の別の例は、摂動に対するシステムの複雑な応答における識別されたファーマコフォアの役割に関する直感を得るために、薬学者などによる専門知識を適用することを含む。

【0549】

新しい分子の設計のためのファーマコフォアを予測するためのモデルを実施し、モデルは、Teverisky類似性を使用してスコアに基づいて選択された介入ライブラリからの低分子の特徴化を含み、ファーマコフォアが化学構造に含まれているかどうかを示す表

10

20

30

40

50

現を達成した。この表現（化学フィンガープリント）を、実施例 1 のモジュール 7 8 のモデル 6 0 1 に入力した。実施例 1 で識別された脂肪標的化ファーマコフォアを使用して、実施例 1 のモジュール 7 8 についてのモデルを使用して、既知のピペリジン含有化合物（「K P C C」）の脂肪標的化ファーマコフォアの関連性を決定し、分離して、0 . 0 4 0 6 4 ~ 0 . 0 4 6 3 3 の範囲の活性化スコアを有する脂肪酸モジュールの転写活性化を観察した。

#### 【0 5 5 0】

実施例 4 . 潜在的な細胞挙動に基づく合成ヒット化合物の生成。

試験事例として、本明細書において「6 つの合成ヒット」と称される新たに合成された低分子ヒットのうち 6 つは、インビトロ及びインビボで検証された脂肪細胞ベージング化合物及びその潜在空間表現に基づいて設計された。6 つの合成ヒットの各々は、ヒト前脂肪細胞上の所望の細胞挙動変化を誘発した。まず、K P C C クラスターのファーマコフォアを識別した。次いで、分子を、新規の生物学的等価体の組み込みとともに、このクラスターにおけるファーマコフォアの濃縮によって設計し、これが 6 つの合成ヒットの最終的な設計につながった。これらの 6 つの構造的に多様な合成ヒットの目標は、K P C C を含む既知の化学物質（K C E）と同じ細胞挙動効果を誘導することであった。図 1 3 の概略図に示されるように、細胞挙動効果は、ヒト前脂肪細胞を 1  $\mu$  M の K P C C 及び 6 つの合成ヒットで 2 4 時間処理し、s c R N A - s e q を使用して遺伝子発現を測定し、上記実施例 1 に記載される脂肪代謝遺伝子モジュールの変化によって発現される細胞応答を評価することによって決定した（モジュール 7 8）。例えば、脂肪代謝モジュールにおける遺伝子は、とりわけ、F A B P 3、F D P S、及び L P I N 1 を含む。

10

20

#### 【0 5 5 1】

前脂肪細胞に対するこれらの化合物の影響の評価は、各々の合成ヒットが、K P C C と同じ脂肪代謝遺伝子モジュールを活性化したことを明らかにした（図 1 3；モジュール 7 8 は、ボックス 1 3 0 2 で強調されている）。すなわち、強調表示されたボックス 1 3 0 2 は、図 1 3 のグラフの Y 軸に列挙されたモデルに化合物のフィンガープリントを入力すると、実施例 1 のモジュール 7 8 についてのモデルによって出力される活性化スコアを示す。これらの結果は、所望の細胞挙動を予測可能に標的とするモデルプラットフォームに基づいて合成ヒットを生成する能力に高い信頼性を提供する。特に、本開示のモデル 6 0 1（例えば、実施例 1 のモジュール 7 8 についてのモデル 6 0 1）を使用して、ハイスループットスクリーニング、分子標的に基づく識別若しくは最適化、又は検証のための数百又は数千の新規化合物の合成を必要とせずに、生理学的状態と関連する遺伝子モジュールを標的とする合成ヒットを予測することができる。

30

#### 【0 5 5 2】

引用文献及び代替の実施形態

本明細書で引用される全ての参考文献は、個々の刊行物又は特許又は特許出願の各々が、全ての目的において参照によりその全体が組み込まれることが具体的かつ個々に示されたのと同様に、全ての目的において参照によりその全体が本明細書に組み込まれる。

#### 【0 5 5 3】

本発明は、非一時的なコンピュータ可読記憶媒体に埋め込まれたコンピュータプログラム機構を含むコンピュータプログラム製品として実装され得る。例えば、コンピュータプログラム製品は、図 1 ~ 図 3 及び図 7 ~ 図 9 の任意の組み合わせで示されるプログラムモジュールを含み得る。これらのプログラムモジュールは、C D - R O M、D V D、磁気ディスクストレージ製品、又は任意の他の非一時的コンピュータ可読データ若しくはプログラムストレージ製品に格納され得る。

40

#### 【0 5 5 4】

当業者には明らかなように、本発明の多くの修正及び変形を、その趣旨及び範囲から逸脱することなく行うことができる。本明細書に記載される特定の実施形態は、例としてのみ提供される。実施形態は、本発明の原理及びその実際の用途を最もよく説明するために選択及び説明され、それによって、当業者が、本発明及び企図される特定の使用に適した

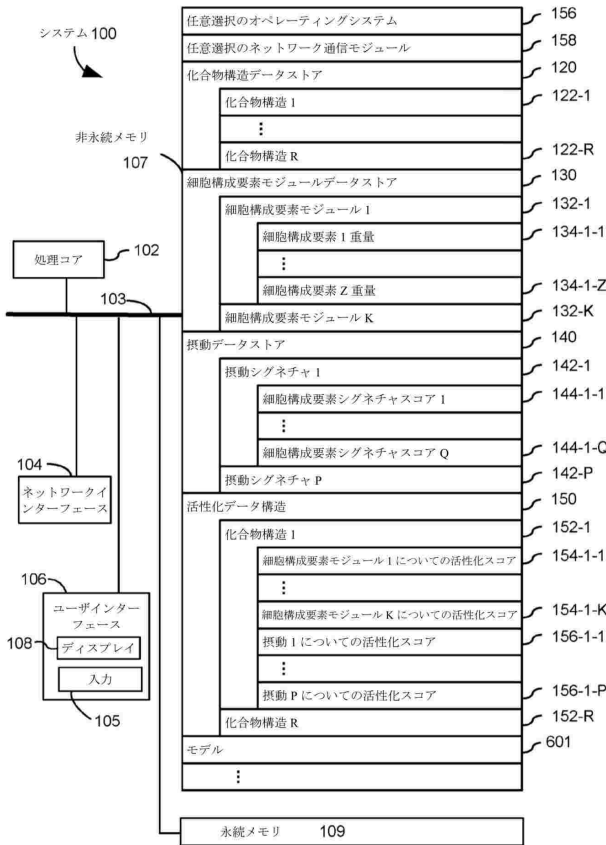
50



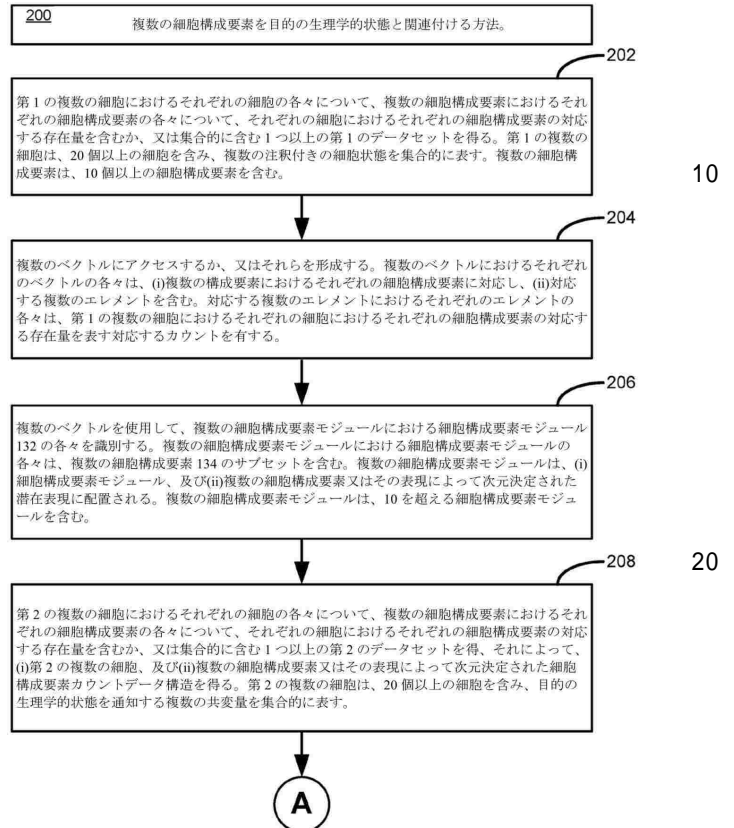
様々な修正を伴う様々な実施形態を最もよく利用できるようにする。本発明は、添付の特許請求の範囲の用語、及びそのような特許請求の範囲が権利を有する等価物の全範囲によってのみ限定される。

【図面】

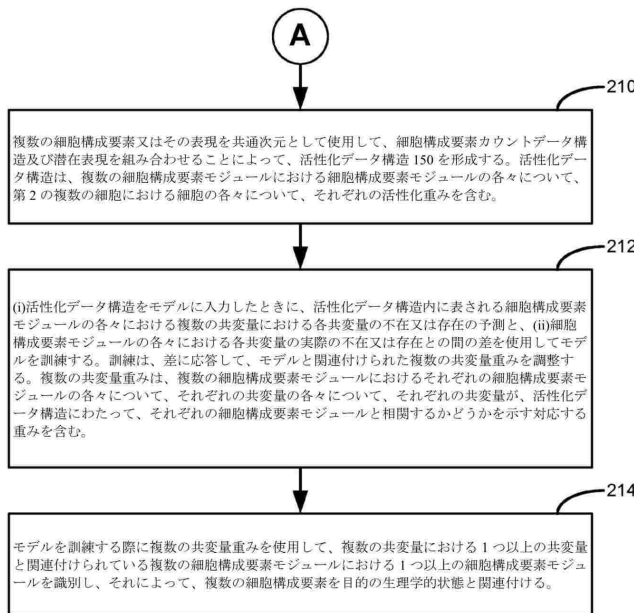
【図 1】



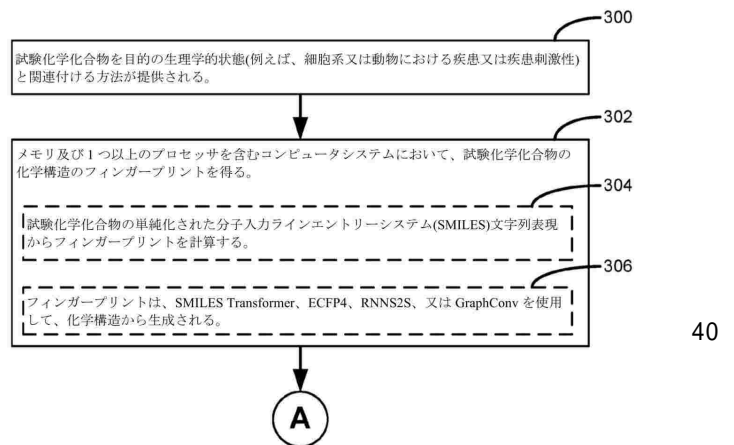
【図 2 A】



【図 2 B】



【図 3 A】

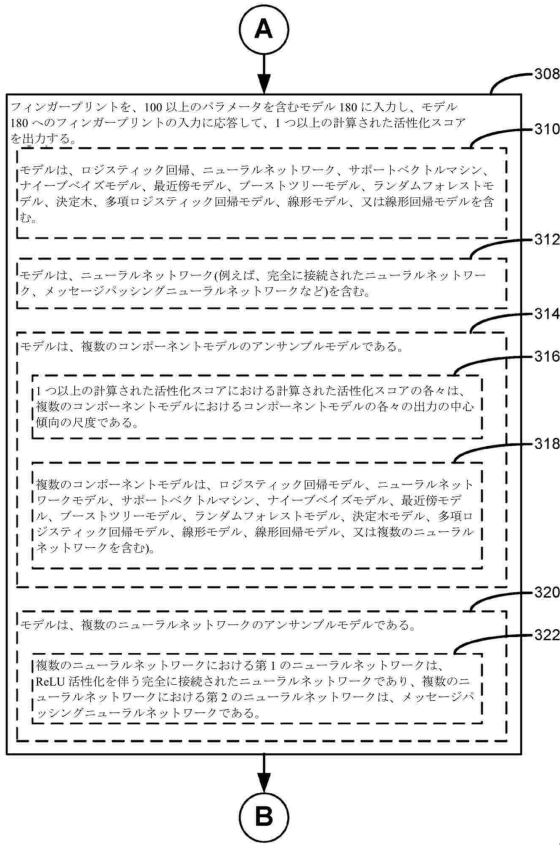


30

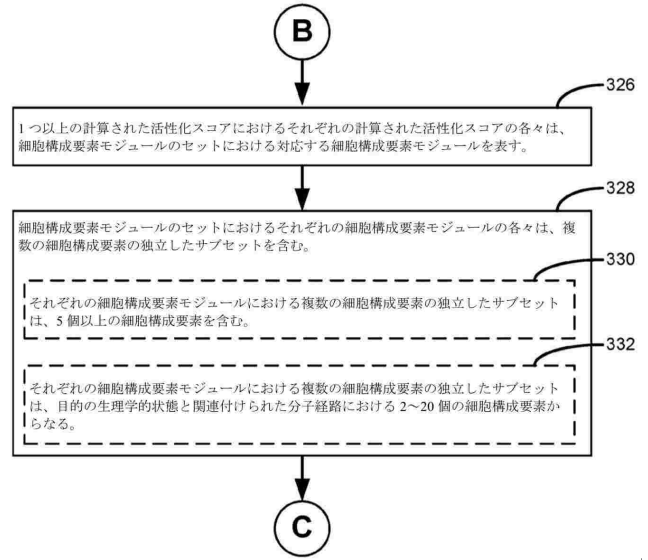
40

50

【図 3 B】



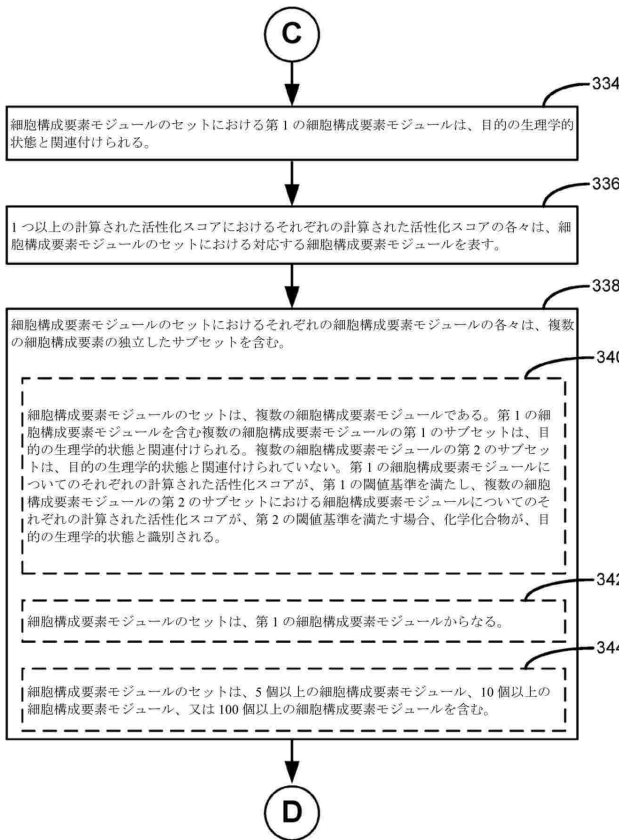
【図 3 C】



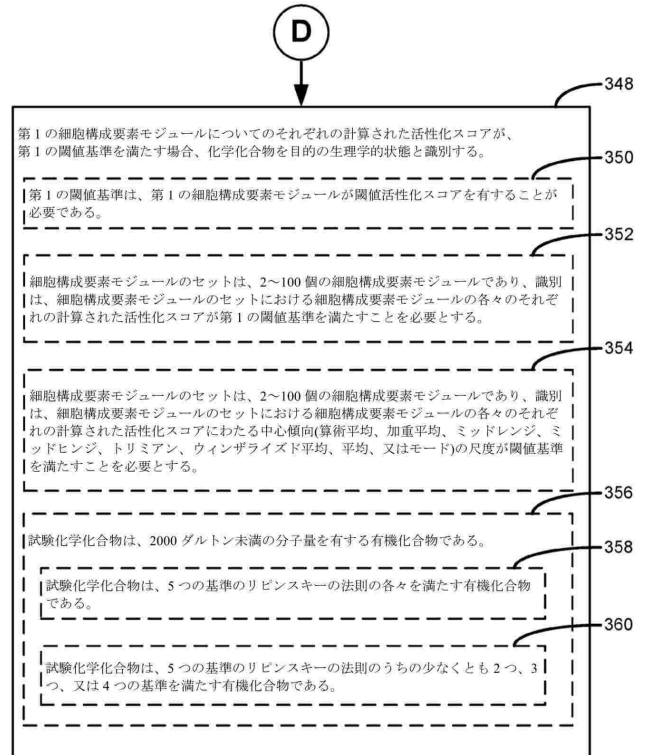
10

20

【図 3 D】



【図 3 E】

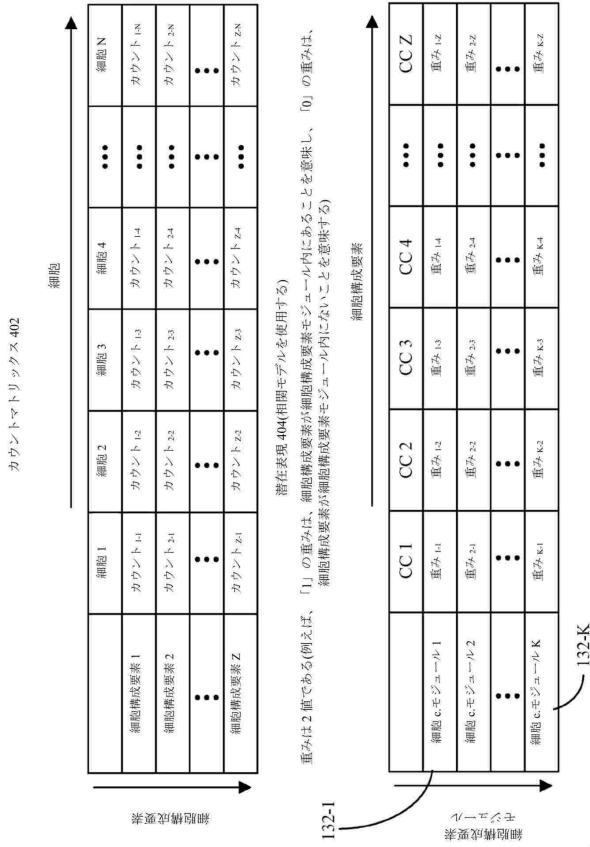


30

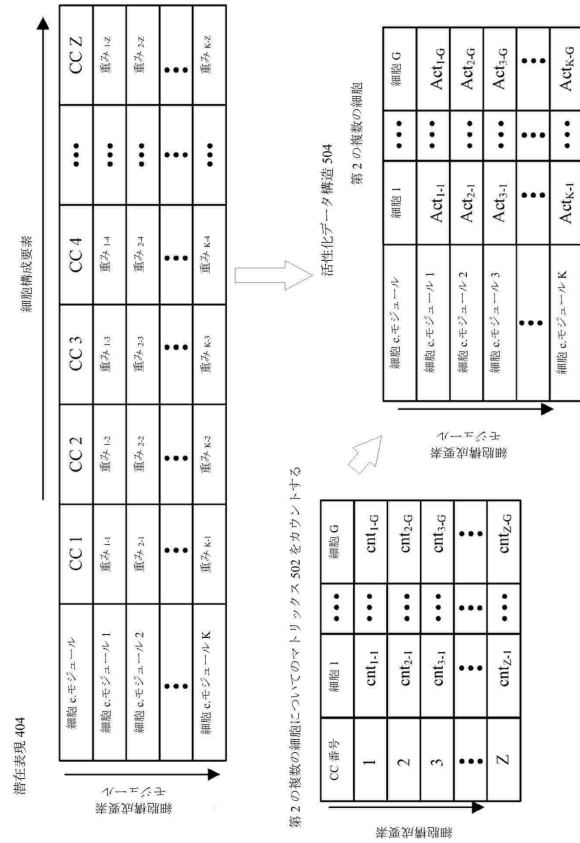
40

50

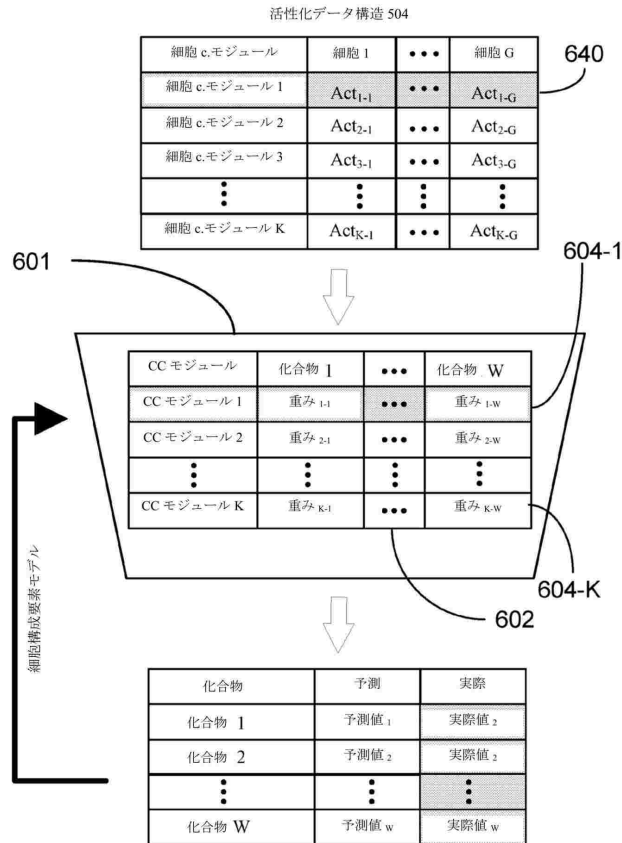
【 図 4 】



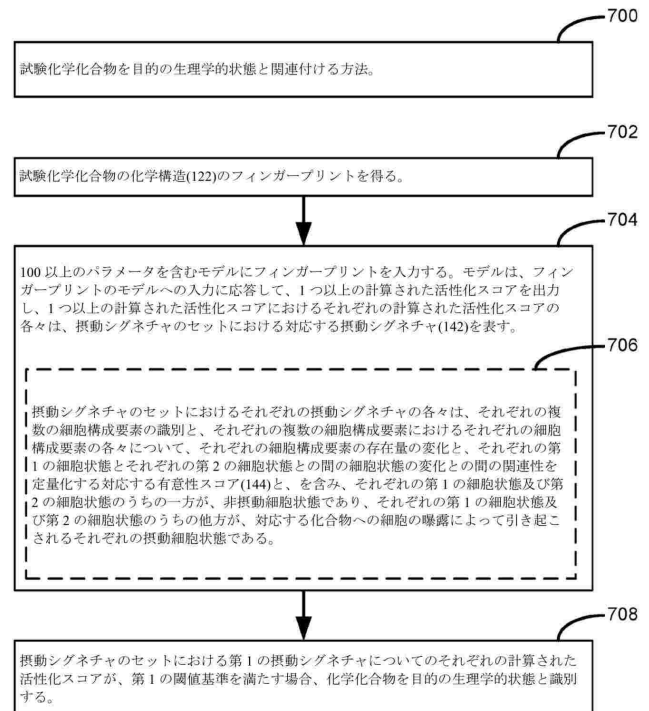
【 図 5 】



【 図 6 】



【 図 7 】



10

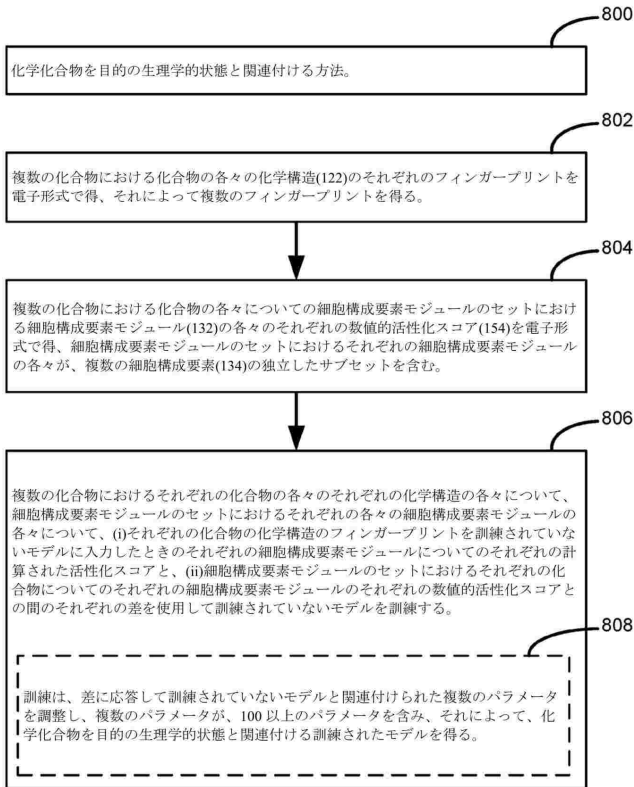
20

30

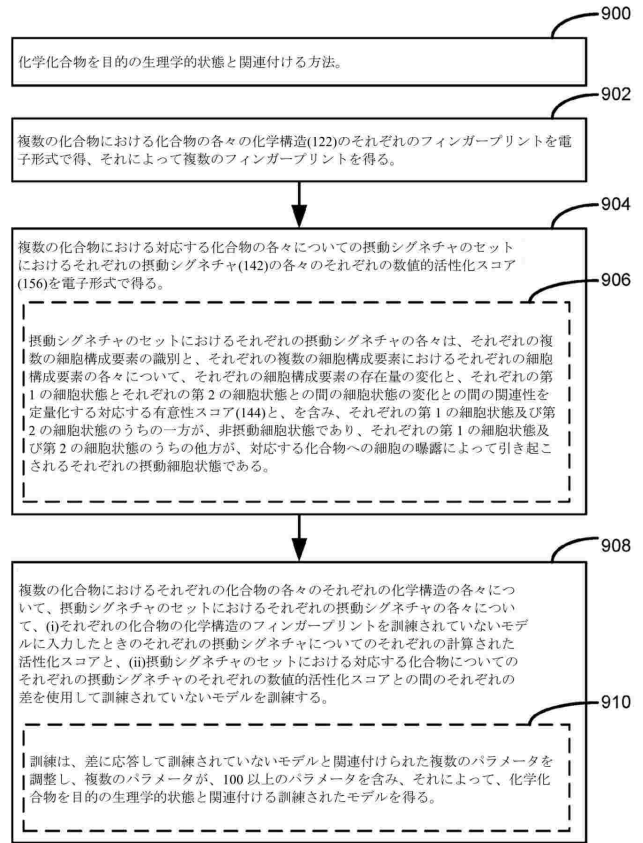
40

50

【 図 8 】



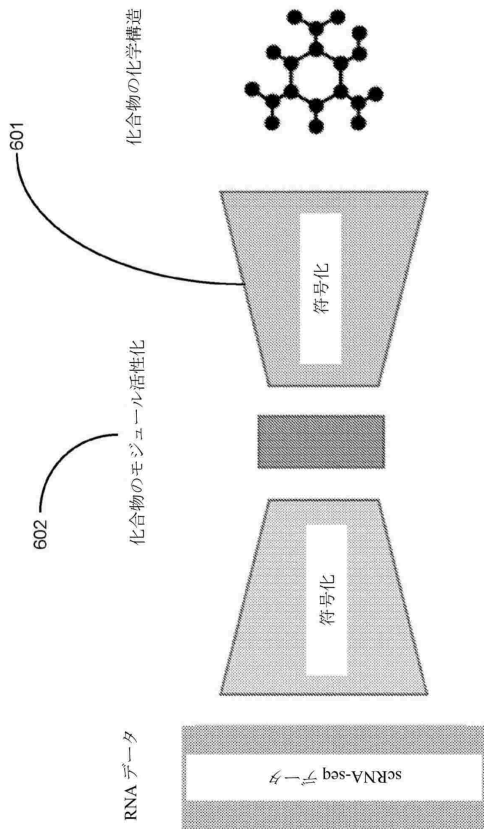
【 図 9 】



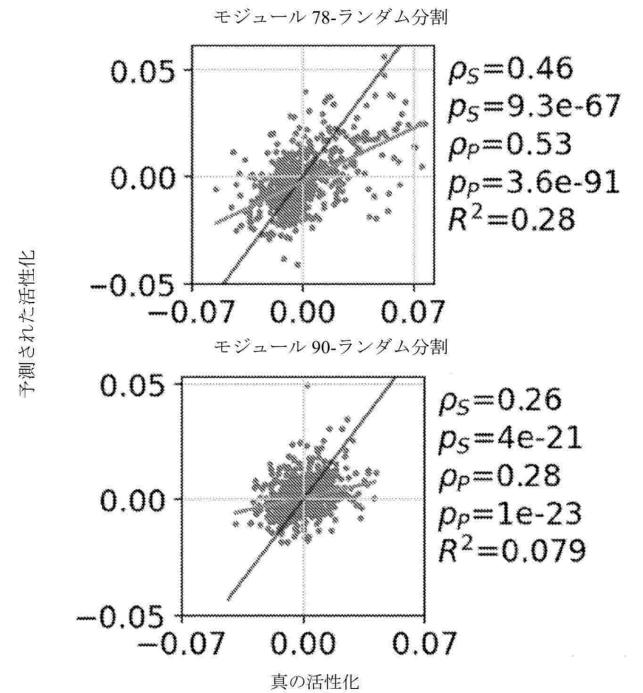
10

20

【 図 10 A 】



【 図 10 B 】

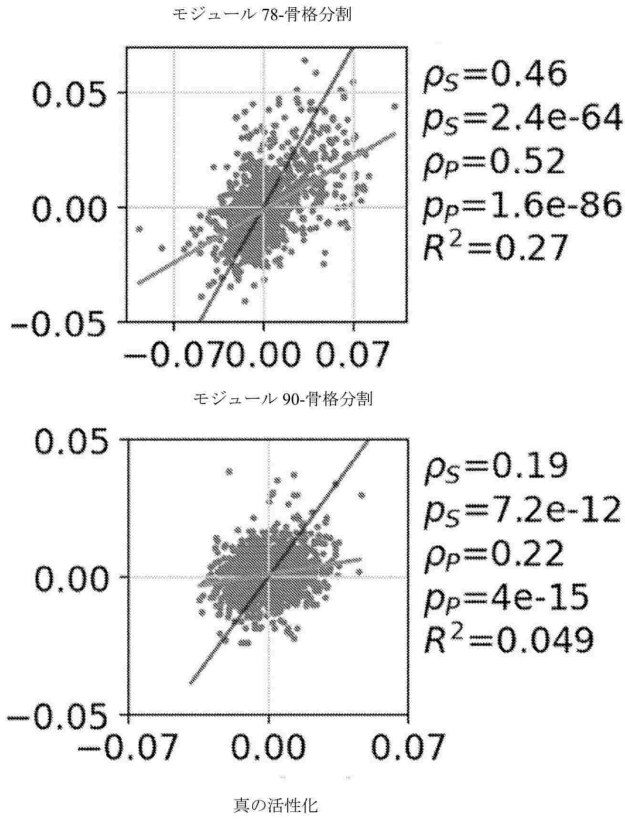


30

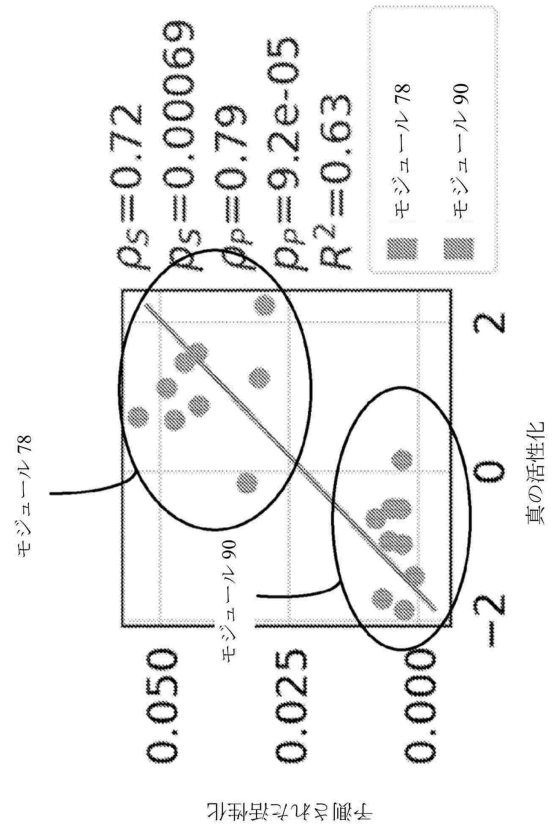
40

50

【 図 1 0 C 】



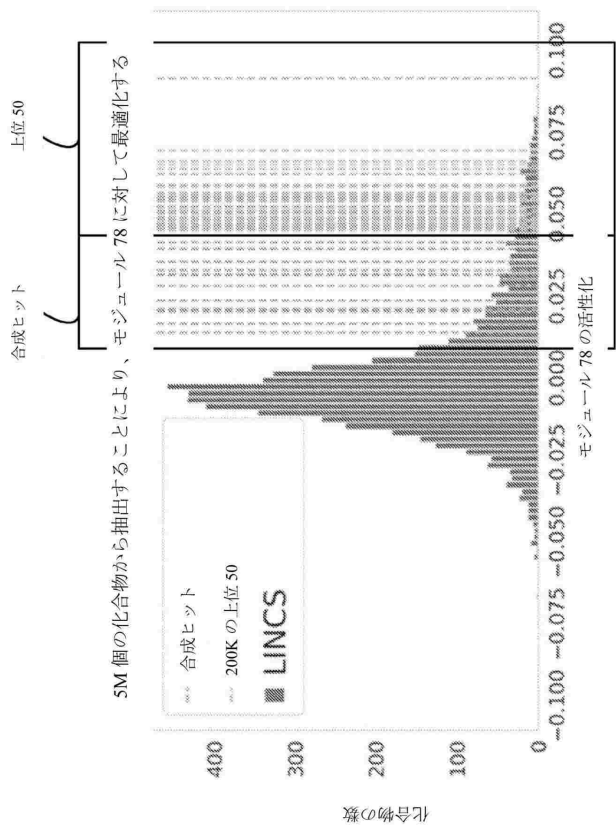
【 図 1 0 D 】



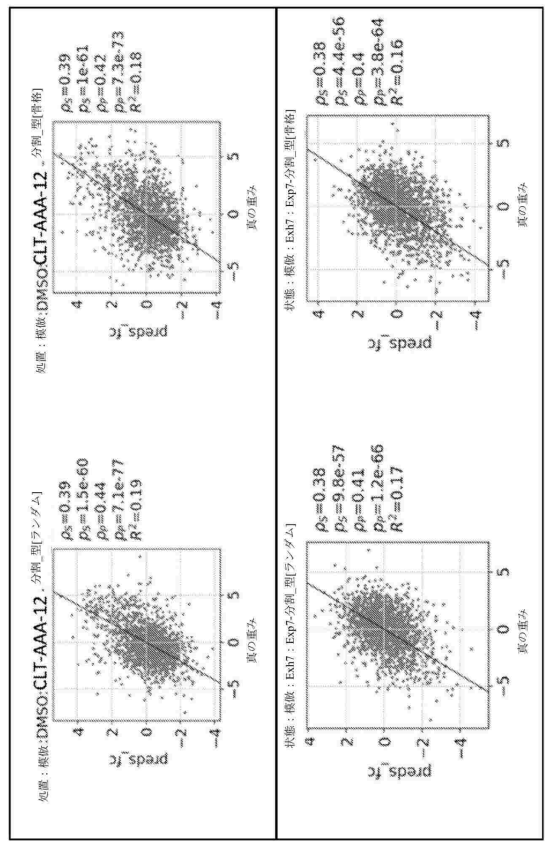
10

20

【 図 1 0 E 】



【 図 1 1 】

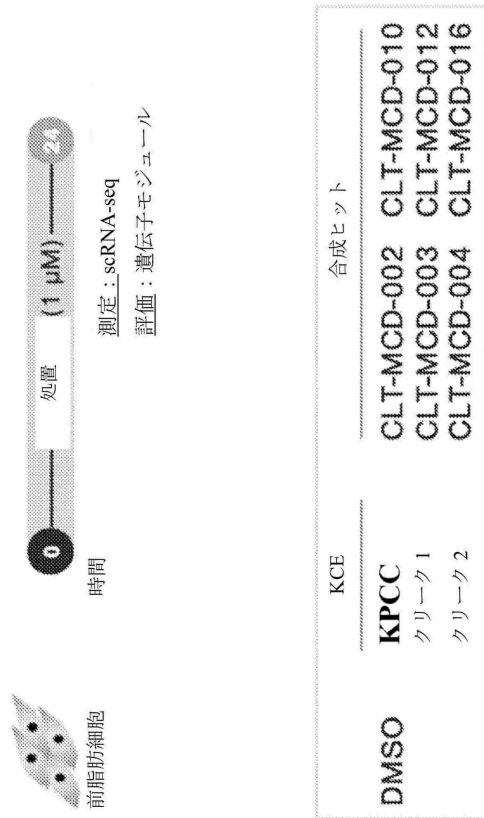


30

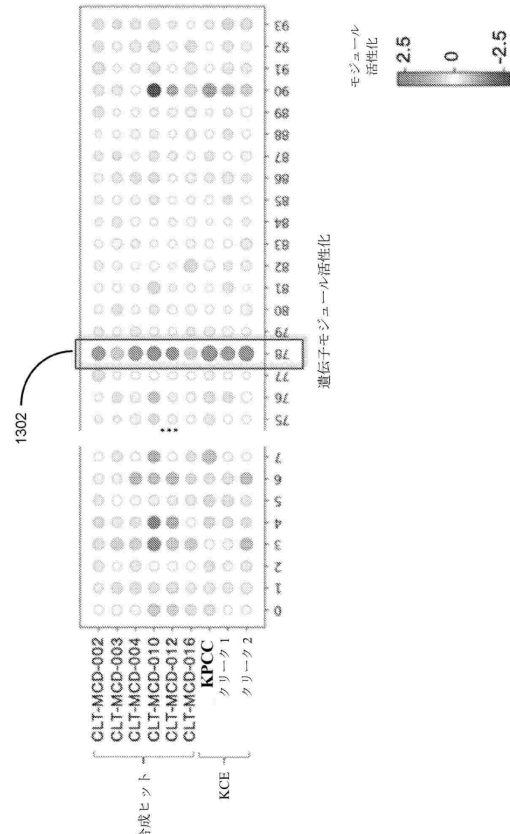
40

50

【 図 1 2 】



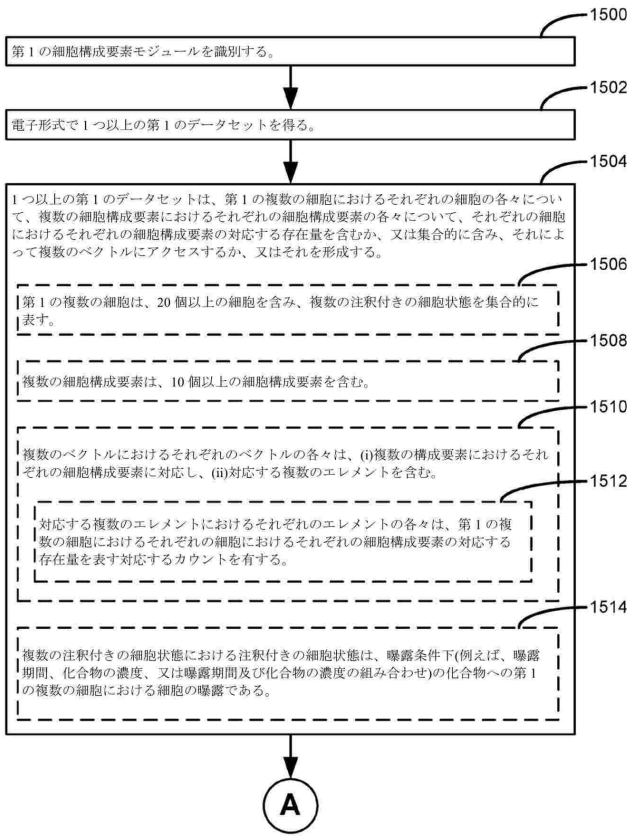
【 図 1 3 】



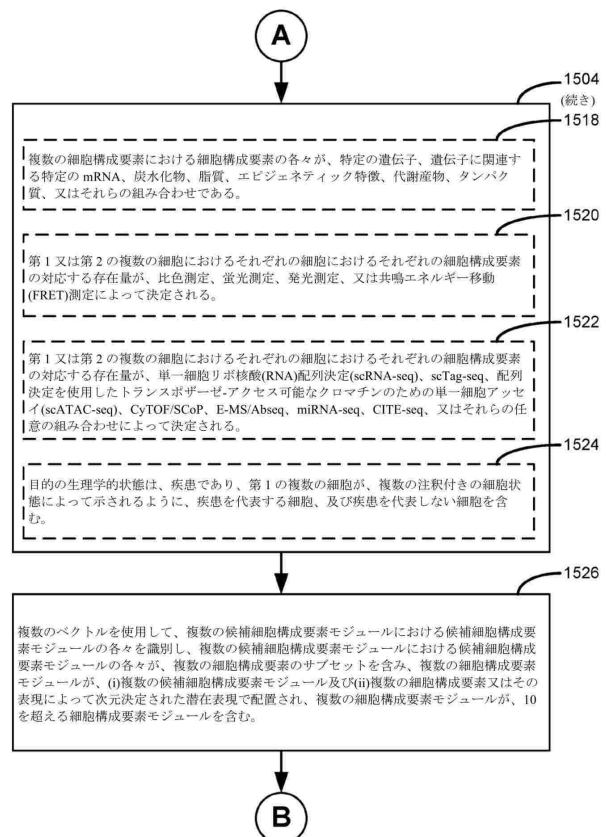
10

20

【 図 1 4 A 】



【 図 1 4 B 】

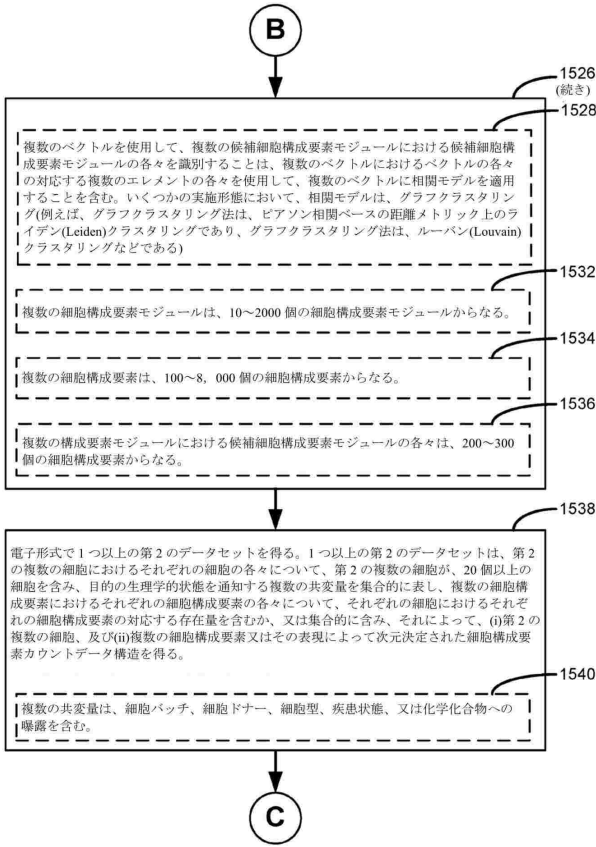


30

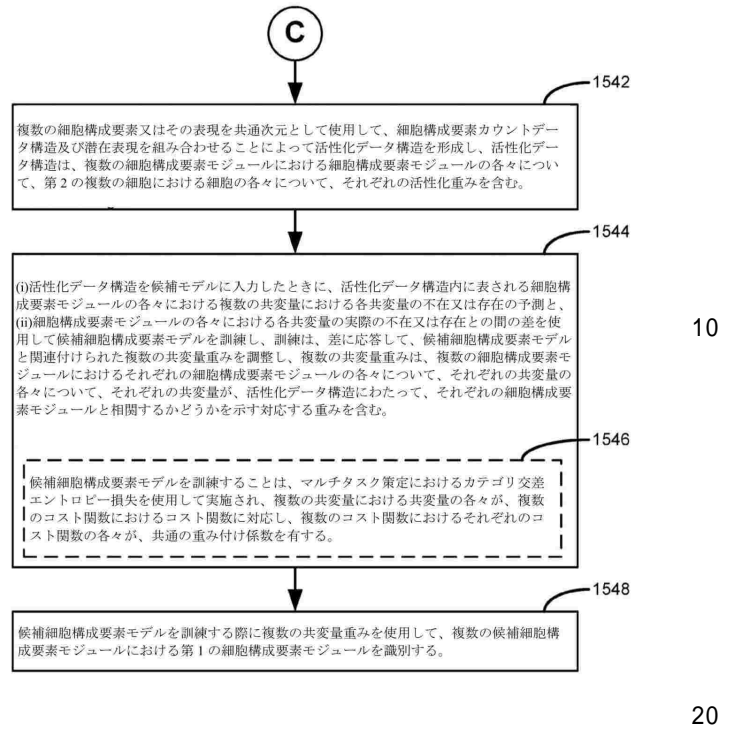
40

50

【 図 1 4 C 】



【 図 1 4 D 】



【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No  
PCT/US2022/033685

<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
INV.	G16B15/30	G16B5/00 G16B25/10 G16B40/20
ADD.		
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, WPI Data		
<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	Lotfollahi Mohammad ET AL: "Learning interpretable cellular responses to complex perturbations in high-throughput screens", bioRxiv, 15 April 2021 (2021-04-15), XP055962485, DOI: 10.1101/2021.04.14.439903 Retrieved from the Internet: URL:https://www.biorxiv.org/content/10.1101/2021.04.14.439903v1.full.pdf [retrieved on 2022-09-19] the whole document in particular "results" and "methods" sections; figures 1-3 ----- -/--	1-159
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" 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	"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	
"E" earlier application or patent but published on or after the international filing date	"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	
"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)	"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	
"O" document referring to an oral disclosure, use, exhibition or other means	"&" document member of the same patent family	
"P" document published prior to the international filing date but later than the priority date claimed		
Date of the actual completion of the international search <b>20 September 2022</b>	Date of mailing of the international search report <b>28/09/2022</b>	
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 <b>Rákossy, Z</b>	

10

20

30

40

1

50



## INTERNATIONAL SEARCH REPORT

International application No PCT/US2022/033685
---

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p><b>DURAN-FRIGOLA MIQUEL ET AL:</b> "Extending the small-molecule similarity principle to all levels of biology with the Chemical Checker", NATURE BIOTECHNOLOGY, NATURE PUBLISHING GROUP US, NEW YORK, vol. 38, no. 9, 18 May 2020 (2020-05-18), pages 1087-1096, XP037524011, ISSN: 1087-0156, DOI: 10.1038/S41587-020-0502-7 [retrieved on 2020-05-18] the whole document in particular "methods" section; page 1093 - page 1095; figures 1, 2, 5; table 1</p> <p style="text-align: center;">-----</p>	1-159
X	<p><b>JO JEONGHEE ET AL:</b> "The message passing neural networks for chemical property prediction on SMILES", METHODS, ACADEMIC PRESS, NL, vol. 179, 21 May 2020 (2020-05-21), pages 65-72, XP086236801, ISSN: 1046-2023, DOI: 10.1016/J.YMETH.2020.05.009 [retrieved on 2020-05-21] the whole document in particular sections 2 and 3; figure 1</p> <p style="text-align: center;">-----</p>	1-159
X	<p><b>US 2020/020419 A1 (KAHVEJIAN AVAK [US] ET AL)</b> 16 January 2020 (2020-01-16) the whole document in particular paragraph [0115] - paragraph [0226]; figures 2, 3</p> <p style="text-align: center;">-----</p>	1-159

1

Form PCT/ISA/210 (continuation of second sheet) (April 2005)

page 2 of 2

10

20

30

40

50

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No <b>PCT/US2022/033685</b>
--

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
<b>US 2020020419 A1</b>	<b>16-01-2020</b>	<b>CA 3103677 A1</b>	<b>23-01-2020</b>
		<b>CN 112424866 A</b>	<b>26-02-2021</b>
		<b>EP 3824080 A1</b>	<b>26-05-2021</b>
		<b>JP 2022501011 A</b>	<b>06-01-2022</b>
		<b>KR 20210031708 A</b>	<b>22-03-2021</b>
		<b>US 2020020419 A1</b>	<b>16-01-2020</b>
		<b>WO 2020018519 A1</b>	<b>23-01-2020</b>

---

10

20

30

40

50

## フロントページの続き

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,N  
E,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,IQ,IR,IS,IT,JM,J  
O,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,M  
Z,NA,NG,NI,NO,NZ,OM,PA,PE,PG,PH,PL,PT,QA,RO,RS,RU,RW,SA,SC,SD,SE,SG,SK,SL,ST,SV,SY,TH,  
TJ,TM,TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

(特許庁注：以下のものは登録商標)

1 . B L U E T O O T H

2 . A N D R O I D

3 . L i n u x

4 . U N I X

5 . O S X

6 . W I N D O W S

7 . V X W O R K S

8 . T E N S O R F L O W

9 . W C D M A

アメリカ合衆国，マサチューセッツ州 0 2 1 4 2 ，ケンブリッジ，ケンブリッジ パークウェイ  
5 5 エイス フロア フラッグシップ パイオニアリング イノベーションズ シックス，エルエルシー

(72)発明者 ハダド，レイジー

アメリカ合衆国，マサチューセッツ州 0 2 1 4 2 ，ケンブリッジ，ケンブリッジ パークウェイ  
5 5 エイス フロア フラッグシップ パイオニアリング イノベーションズ シックス，エルエルシー

(72)発明者 プラギス，ニコラス マッカートニー

アメリカ合衆国，マサチューセッツ州 0 2 1 4 2 ，ケンブリッジ，ケンブリッジ パークウェイ  
5 5 エイス フロア フラッグシップ パイオニアリング イノベーションズ シックス，エルエルシー

F ターム ( 参考 ) 4B063 QA01 QA13 QQ08 QQ53 QR08 QR42 QR62 QS34 QX02  
4B065 AA93X CA46