

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2020-525892  
(P2020-525892A)

(43) 公表日 令和2年8月27日(2020.8.27)

(51) Int.Cl. F I テーマコード(参考)  
**G 1 6 B 20/00 (2019.01)** G 1 6 B 20/00  
**G 0 6 N 3/08 (2006.01)** G 0 6 N 3/08

審査請求有 予備審査請求有 (全 157 頁)

(21) 出願番号 特願2019-567719 (P2019-567719)  
 (86) (22) 出願日 平成30年10月15日(2018.10.15)  
 (85) 翻訳文提出日 令和1年12月25日(2019.12.25)  
 (86) 国際出願番号 PCT/US2018/055840  
 (87) 国際公開番号 W02019/079166  
 (87) 国際公開日 平成31年4月25日(2019.4.25)  
 (31) 優先権主張番号 62/573,144  
 (32) 優先日 平成29年10月16日(2017.10.16)  
 (33) 優先権主張国・地域又は機関 米国(US)  
 (31) 優先権主張番号 62/573,149  
 (32) 優先日 平成29年10月16日(2017.10.16)  
 (33) 優先権主張国・地域又は機関 米国(US)

(71) 出願人 500358711  
 イルミナ インコーポレイテッド  
 アメリカ合衆国 カリフォルニア州 92  
 122 サンディエゴ イルミナ ウエイ  
 5200  
 (74) 代理人 100108453  
 弁理士 村山 靖彦  
 (74) 代理人 100110364  
 弁理士 実広 信哉  
 (74) 代理人 100133400  
 弁理士 阿部 達彦  
 (72) 発明者 ホン・ガオ  
 アメリカ合衆国・カリフォルニア・921  
 22・サン・ディエゴ・イルミナ・ウエイ  
 ・5200

最終頁に続く

(54) 【発明の名称】 深層畳み込みニューラルネットワークを訓練するための深層学習ベースの技法

(57) 【要約】

開示される技術は、バリエーション分類のための畳み込みニューラルネットワークベースの分類器を構築することに関する。具体的には、開示される技術は、畳み込みニューラルネットワークベースの分類器の出力を対応するグラウンドトゥールズラベルと漸進的に照合する逆伝播ベースの勾配更新技法を使用して、訓練データについて畳み込みニューラルネットワークベースの分類器を訓練することに関する。畳み込みニューラルネットワークベースの分類器は、残差ブロックのグループを備え、残差ブロックの各グループは、残差ブロックの中の畳み込みフィルタの数、残差ブロックの畳み込みウィンドウサイズ、および残差ブロックの膨張畳み込み率によってパラメータ化され、畳み込みウィンドウのサイズは残差ブロックのグループ間で変動し、膨張畳み込み率は残差ブロックのグループ間で変動する。訓練データは、良性バリエーションおよび病原性バリエーションから生成された、翻訳された配列ペアの良性訓練例および病原性訓練例を含む。

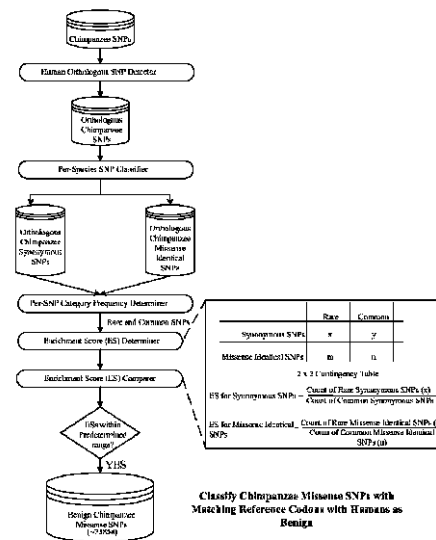


FIG. 46

**【特許請求の範囲】****【請求項1】**

バリエーション病原性分類器を構築する方法であって、

基準タンパク質配列および代替タンパク質配列の良性訓練例のペアと病原性訓練例のペアとを入力として使用して、メモリに結合された多数のプロセッサ上で実行される、畳み込みニューラルネットワークベースのバリエーション病原性分類器を訓練するステップを含み、前記代替タンパク質配列が良性バリエーションおよび病原性バリエーションから生成され、

前記良性バリエーションが、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準コドン配列を共有する代替的なヒト以外の霊長類コドン配列上で発生するヒト以外の霊長類ミスセンスバリエーションとを含む、方法。

10

**【請求項2】**

前記一般的なヒトミスセンスバリエーションが、少なくとも100000人のヒトからサンプリングされたヒト集団バリエーションデータセットにわたって0.1%より高いマイナーアレル頻度(MAFと省略される)を有する、請求項1に記載の方法。

**【請求項3】**

前記サンプリングされたヒトが異なるヒト亜集団に属し、前記一般的なヒトミスセンスバリエーションがそれぞれのヒト亜集団バリエーションデータセット内で0.1%より高いMAFを有する、請求項2に記載の方法。

**【請求項4】**

前記ヒト亜集団が、アフリカ人/アフリカ系アメリカ人(AFRと省略される)、アメリカ人(AMRと省略される)、アシュケナーズ系ユダヤ人(ASJと省略される)、東アジア人(EASと省略される)、フィンランド人(FINと省略される)、フィンランド人以外のヨーロッパ人(NFEと省略される)、南アジア人(SASと省略される)、および他(OTHと省略される)を含む、請求項3に記載の方法。

20

**【請求項5】**

前記ヒト以外の霊長類ミスセンスバリエーションが、チンパンジー、ボノボ、ゴリラ、B.オランウータン、S.オランウータン、アカゲザル、およびマーモセットを含む、複数のヒト以外の霊長類の種からのミスセンスバリエーションを含む、請求項1に記載の方法。

**【請求項6】**

エンリッチメント分析に基づいて、前記良性バリエーションに特定のヒト以外の霊長類の種のミスセンスバリエーションを含めるために、前記特定のヒト以外の霊長類の種を受け入れるステップをさらに含み、前記エンリッチメント分析が、前記特定のヒト以外の霊長類の種に対して、前記特定のヒト以外の霊長類の種の同義バリエーションの第1のエンリッチメントスコアを前記特定のヒト以外の霊長類の種のミスセンス同一バリエーションの第2のエンリッチメントスコアと比較することを含み、

30

ミスセンス同一バリエーションが、ヒトと一致する基準コドン配列および代替コドン配列を共有するミスセンスバリエーションであり、

前記第1のエンリッチメントスコアが、0.1%より高いMAFを伴う一般的な同義バリエーションに対する0.1%より低いMAFを伴う稀な同義バリエーションの比を決定することによって作り出され、

40

前記第2のエンリッチメントスコアが、0.1%より高いMAFを伴う一般的なミスセンス同一バリエーションに対する0.1%より低いMAFを伴う稀なミスセンス同一バリエーションの比を決定することによって作り出される、請求項1に記載の方法。

**【請求項7】**

稀な同義バリエーションがシングルトンバリエーションを含む、請求項6に記載の方法。

**【請求項8】**

前記第1のエンリッチメントスコアと前記第2のエンリッチメントスコアとの差が所定の範囲内にあり、前記良性バリエーションに前記特定のヒト以外の霊長類のミスセンスバリエーションを含めるために、前記特定のヒト以外の霊長類の種を受け入れるステップをさらに含み、請求項6に記載の方法。

50

## 【請求項 9】

前記差が前記所定の範囲にあることが、前記ミスセンス同一バリエーションが前記同義バリエーションと同じ程度の自然選択を受けており、したがって前記同義バリエーションと同じくらい良性であることを示す、請求項6に記載の方法。

## 【請求項 10】

前記エンリッチメント分析を繰り返し適用して、前記良性バリエーションに前記ヒト以外の霊長類の種のミスセンスバリエーションを含めるために複数のヒト以外の霊長類の種を受け入れるステップをさらに含む、請求項6に記載の方法。

## 【請求項 11】

前記ヒト以外の霊長類の種の各々に対する同義バリエーションの第1のエンリッチメントスコアとミスセンス同一バリエーションの第2のエンリッチメントスコアを比較するために、相同性のカイ二乗検定を使用するステップをさらに含む、請求項1に記載の方法。

10

## 【請求項 12】

前記ヒト以外の霊長類ミスセンスバリエーションのカウントが少なくとも100000である、請求項1に記載の方法。

## 【請求項 13】

前記ヒト以外の霊長類ミスセンスバリエーションの前記カウントが385236である、請求項12に記載の方法。

## 【請求項 14】

前記一般的なヒトミスセンスバリエーションのカウントが少なくとも50000である、請求項1に記載の方法。

20

## 【請求項 15】

前記一般的なヒトミスセンスバリエーションの前記カウントが83546である、請求項14に記載の方法。

## 【請求項 16】

バリエーション病原性分類器を構築するためのコンピュータプログラム命令が焼かれた非一時的コンピュータ可読記憶媒体であって、プロセッサで実行されると、

基準タンパク質配列および代替タンパク質配列の良性訓練例のペアと病原性訓練例のペアとを入力として使用して、メモリに結合された多数のプロセッサ上で実行される、畳み込みニューラルネットワークベースのバリエーション病原性分類器を訓練するステップを含む方法を実施し、前記代替タンパク質配列が良性バリエーションおよび病原性バリエーションから生成され、

30

前記良性バリエーションが、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準コドン配列を共有する代替的なヒト以外の霊長類コドン配列上で発生するヒト以外の霊長類ミスセンスバリエーションとを含む、非一時的コンピュータ可読記憶媒体。

## 【請求項 17】

エンリッチメント分析に基づいて、前記良性バリエーションに特定のヒト以外の霊長類の種のミスセンスバリエーションを含めるために、前記特定のヒト以外の霊長類の種を受け入れることをさらに含む前記方法を実施し、前記エンリッチメント分析が、前記特定のヒト以外の霊長類の種に対して、前記特定のヒト以外の霊長類の種の同義バリエーションの第1のエンリッチメントスコアを前記特定のヒト以外の霊長類の種のミスセンス同一バリエーションの第2のエンリッチメントスコアと比較することを含み、

40

ミスセンス同一バリエーションが、ヒトと一致する基準コドン配列および代替コドン配列を共有するミスセンスバリエーションであり、

前記第1のエンリッチメントスコアが、0.1%より高いMAFを伴う一般的な同義バリエーションに対する0.1%より低いMAFを伴う稀な同義バリエーションの比を決定することによって作り出され、

前記第2のエンリッチメントスコアが、0.1%より高いMAFを伴う一般的なミスセンス同一バリエーションに対する0.1%より低いMAFを伴う稀なミスセンス同一バリエーションの比を決定することによって作り出される、請求項16に記載の非一時的コンピュータ可読記憶媒体。

50

## 【請求項18】

前記ヒト以外の霊長類の種の各々に対する同義バリエーションの第1のエンリッチメントスコアとミスセンス同一バリエーションの第2のエンリッチメントスコアを比較するために、相同性のカイ二乗検定を使用することをさらに含む前記方法を実施する、請求項16に記載の非一時的コンピュータ可読記憶媒体。

## 【請求項19】

バリエーション病原性分類器を構築するためのコンピュータ命令がロードされたメモリに結合された1つまたは複数のプロセッサを含むシステムであって、前記プロセッサで実行されると、前記命令が、

基準タンパク質配列および代替タンパク質配列の良性訓練例のペアと病原性訓練例のペアとを入力として使用して、メモリに結合された多数のプロセッサ上で実行される、畳み込みニューラルネットワークベースのバリエーション病原性分類器を訓練することを含む活動を実施し、前記代替タンパク質配列が良性バリエーションおよび病原性バリエーションから生成され、

前記良性バリエーションが、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準コドン配列を共有する代替的なヒト以外の霊長類コドン配列上で発生するヒト以外の霊長類ミスセンスバリエーションとを含む、システム。

## 【請求項20】

エンリッチメント分析に基づいて、前記良性バリエーションに特定のヒト以外の霊長類の種のミスセンスバリエーションを含めるために、前記特定のヒト以外の霊長類の種を受け入れる活動をさらに実施し、前記エンリッチメント分析が、前記特定のヒト以外の霊長類の種に対して、前記特定のヒト以外の霊長類の種の同義バリエーションの第1のエンリッチメントスコアを前記特定のヒト以外の霊長類の種のミスセンス同一バリエーションの第2のエンリッチメントスコアと比較することを含み、

ミスセンス同一バリエーションが、ヒトと一致する基準コドン配列および代替コドン配列を共有するミスセンスバリエーションであり、

前記第1のエンリッチメントスコアが、0.1%より高いMAFを伴う一般的な同義バリエーションに対する0.1%より低いMAFを伴う稀な同義バリエーションの比を決定することによって作り出され、

前記第2のエンリッチメントスコアが、0.1%より高いMAFを伴う一般的なミスセンス同一バリエーションに対する0.1%より低いMAFを伴う稀なミスセンス同一バリエーションの比を決定することによって作り出される、請求項19に記載のシステム。

## 【請求項21】

バリエーション分類のために畳み込みニューラルネットワークベースの分類器を構築するコンピュータで実施される方法であって、

メモリに結合される多数のプロセッサ上で実行される畳み込みニューラルネットワークベースの分類器を、前記畳み込みニューラルネットワークベースの分類器の出力を対応するグラウンドトゥールズラベルと漸進的に照合する逆伝播ベースの勾配更新技法を使用して、訓練データについて訓練するステップを含み、

前記畳み込みニューラルネットワークベースの分類器が残差ブロックのグループを備え、

残差ブロックの各グループが、前記残差ブロックの中の畳み込みフィルタの数、前記残差ブロックの畳み込みウィンドウサイズ、および前記残差ブロックの膨張畳み込み率によってパラメータ化され、

前記畳み込みウィンドウサイズが残差ブロックの前記グループ間で変動し、

前記膨張畳み込み率が残差ブロックの前記グループ間で変動し、

前記訓練データが、良性バリエーションおよび病原性バリエーションから生成され良性訓練例および病原性訓練例として使用される、翻訳された配列のペアを含み、

前記良性バリエーションが、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準塩基トリプレット配列を共有する代替的なヒト以外の霊長類塩基トリプレット配列上で発生

10

20

30

40

50

するヒト以外の霊長類ミスセンスバリエーションを含む、方法。

【請求項 2 2】

バリエーション分類のために畳み込みニューラルネットワークベースの分類器を構築するためのコンピュータプログラム命令が焼かれた非一時的コンピュータ可読記憶媒体であって、プロセッサで実行されると、前記命令が、

メモリに結合される多数のプロセッサ上で実行される畳み込みニューラルネットワークベースの分類器を、前記畳み込みニューラルネットワークベースの分類器の出力を対応するグラウンドトゥールスラベルと漸進的に照合する逆伝播ベースの勾配更新技法を使用して、訓練データについて訓練するステップを含む方法を実施し、

前記畳み込みニューラルネットワークベースの分類器が残差ブロックのグループを備え

10

、  
残差ブロックの各グループが、前記残差ブロックの中の畳み込みフィルタの数、前記残差ブロックの畳み込みウィンドウサイズ、および前記残差ブロックの膨張畳み込み率によってパラメータ化され、

前記畳み込みウィンドウサイズが残差ブロックの前記グループ間で変動し、

前記膨張畳み込み率が残差ブロックの前記グループ間で変動し、

前記訓練データが、良性バリエーションおよび病原性バリエーションから生成され良性訓練例および病原性訓練例として使用される、翻訳された配列のペアを含み、

前記良性バリエーションが、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準塩基トリプレット配列を共有する代替的なヒト以外の霊長類塩基トリプレット配列上で発生するヒト以外の霊長類ミスセンスバリエーションを含む、非一時的コンピュータ可読記憶媒体。

20

【請求項 2 3】

バリエーション分類のために畳み込みニューラルネットワークベースの分類器を構築するためのコンピュータ命令がロードされたメモリに結合される1つまたは複数のプロセッサを含むシステムであって、前記プロセッサで実行されると、前記命令が、

メモリに結合される多数のプロセッサ上で実行される畳み込みニューラルネットワークベースの分類器を、前記畳み込みニューラルネットワークベースの分類器の出力を対応するグラウンドトゥールスラベルと漸進的に照合する逆伝播ベースの勾配更新技法を使用して、訓練データについて訓練することを含む活動を実施し、

30

前記畳み込みニューラルネットワークベースの分類器が残差ブロックのグループを備え

、  
残差ブロックの各グループが、前記残差ブロックの中の畳み込みフィルタの数、前記残差ブロックの畳み込みウィンドウサイズ、および前記残差ブロックの膨張畳み込み率によってパラメータ化され、

前記畳み込みウィンドウサイズが残差ブロックの前記グループ間で変動し、

前記膨張畳み込み率が残差ブロックの前記グループ間で変動し、

前記訓練データが、良性バリエーションおよび病原性バリエーションから生成され良性訓練例および病原性訓練例として使用される、翻訳された配列のペアを含み、

前記良性バリエーションが、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準塩基トリプレット配列を共有する代替的なヒト以外の霊長類塩基トリプレット配列上で発生するヒト以外の霊長類ミスセンスバリエーションを含む、システム。

40

【発明の詳細な説明】

【技術分野】

【0001】

付録

付録には、発明者らが著述した論文に列挙される潜在的な関連する参考文献の目録が含まれる。その論文の主題は、本出願がその優先権を主張する/その利益を主張する米国仮出願において扱われる。これらの参考文献は、要求に応じて訴訟代理人に対して利用可能にされることが可能であり、またはGlobal Dossierを介して入手可能であることがある。

50

その論文は最初の列挙される参考文献である。

【 0 0 0 2 】

優先出願

本出願は、2017年10月16日に出願された、Hong Gao、Kai-How Farh、Laksshman Sundaram、およびJeremy Francis McRaeによる「Training a Deep Pathogenicity Classifier Using Large-Scale Benign Training Data」という表題の米国仮特許出願第62/573,144号(代理人整理番号第ILLM 1000-1/IP-1611-PRV)、2017年10月16日に出願された、Kai-How Farh、Laksshman Sundaram、Samskruthi Reddy Padigepati、およびJeremy Francis McRaeによる「Pathogenicity Classifier Based On Deep Convolutional Neural Networks (CNNs)」という表題の米国仮特許出願第62/573,149号(代理人整理番号第ILLM 1000-2/IP-1612-PRV)、2017年10月16日に出願された、Hong Gao、Kai-How Farh、Laksshman Sundaram、およびJeremy Francis McRaeによる「Deep Semi-Supervised Learning that Generates Large-Scale Pathogenic Training Data」という表題の米国仮特許出願第62/573,153号(代理人整理番号第ILLM 1000-3 /IP-1613-PRV)、および、2017年11月7日に出願された、Hong Gao、Kai-How Farh、およびLaksshman Sundaramによる「Pathogenicity Classification of Genomic Data Using Deep Convolutional Neural Networks (CNNs)」という表題の米国仮特許出願第62/582,898号(代理人整理番号第ILLM 1000-4/IP-1618-PRV)の優先権または利益を主張する。これらの仮出願は、すべての目的のために本明細書において参照により引用される。

10

【 0 0 0 3 】

引用

以下は、本明細書に完全に記載されるかのようにすべての目的のために参照により引用される。

20

【 0 0 0 4 】

後にPCT出願公開第WO(未定)号として公開される、2018年10月15日に同時に出願された、Laksshman Sundaram、Kai-How Farh、Hong Gao、Samskruthi Reddy Padigepati、およびJeremy Francis McRaeによる、「DEEP CONVOLUTIONAL NEURAL NETWORKS FOR VARIANT CLASSIFICATION」という表題のPCT特許出願第PCT/US2018/(未定)号(代理人整理番号ILLM 1000-9/IP-1612-PCT)。

【 0 0 0 5 】

後にPCT出願第WO(未定)号として公開される、2018年10月15日に同時に出願された、Laksshman Sundaram、Kai-How Farh、Hong Gao、およびJeremy Francis McRaeによる、「SEMI-SUPERVISED LEARNING FOR TRAINING AN ENSEMBLE OF DEEP CONVOLUTIONAL NEURAL NETWORKS」という表題の国際特許出願第PCT/US18/(未定)号(代理人整理番号第ILLM 1000-10/IP-1613-PCT)。

30

【 0 0 0 6 】

同時に出願された、Hong Gao、Kai-How Farh、Laksshman Sundaram、およびJeremy Francis McRaeによる、「DEEP LEARNING-BASED TECHNIQUES FOR TRAINING DEEP CONVOLUTIONAL NEURAL NETWORKS」という表題の米国非仮特許出願(代理人整理番号ILLM 1000-5/IP-1611-US)。

40

【 0 0 0 7 】

同時に出願された、Laksshman Sundaram、Kai-How Farh、Hong Gao、およびJeremy Francis McRaeによる、「DEEP CONVOLUTIONAL NEURAL NETWORKS FOR VARIANT CLASSIFICATION」という表題の米国非仮特許出願(代理人整理番号ILLM 1000-6/IP-1612-US)。

【 0 0 0 8 】

同時に出願された、Laksshman Sundaram、Kai-How Farh、Hong Gao、およびJeremy Francis McRaeによる、「SEMI-SUPERVISED LEARNING FOR TRAINING AN ENSEMBLE OF DEEP CONVOLUTIONAL NEURAL NETWORKS」という表題の米国非仮特許出願(代理人整理番号ILLM 1000-7/IP-1613-US)。

【 0 0 0 9 】

50

文書1 - A.V.D.Oord, S.Dieleman, H.Zen, K.Simonyan, O.Vinyals, A.Graves, N.Kalchbrenner, A.Senior, およびK.Kavukcuoglu, 「WAVENET: A GENERATIVE MODEL FOR RAW AUDIO」、arXiv:1609.03499、2016

【 0 0 1 0 】

文書2 - S.O.Arik, M.Chrzanowski, A.Coates, G.Diamos, A.Gibiansky, Y.Kang, X.Li, J.Miller, A.Ng, J.Raiman, S.Sengupta, およびM.Shoeybi, 「DEEP VOICE: REAL-TIME NEURAL TEXT-TO-SPEECH」、arXiv:1702.07825、2017

【 0 0 1 1 】

文書3 - F.YuおよびV.Koltun, 「MULTI-SCALE CONTEXT AGGREGATION BY DILATED CONVOLUTIONS」、arXiv:1511.07122、2016

【 0 0 1 2 】

文書4 - K.He, X.Zhang, S.Ren, およびJ.Sun, 「DEEP RESIDUAL LEARNING FOR IMAGE RECOGNITION」、arXiv:1512.03385、2015

【 0 0 1 3 】

文書5 - R.K.Srivastava, K.Greff, およびJ.Schmidhuber, 「HIGHWAY NETWORKS」、arXiv:1505.00387、2015

【 0 0 1 4 】

文書6 - G.Huang, Z.Liu, L.van der Maaten, およびK.Q.Weinberger, 「DENSELY CONNECTED CONVOLUTIONAL NETWORKS」、arXiv:1608.06993、2017

【 0 0 1 5 】

文書7 - C.Szegedy, W.Liu, Y.Jia, P.Sermanet, S.Reed, D.Anguelov, D.Erhan, V.Vanhoucke, およびA.Rabinovich, 「GOING DEEPER WITH CONVOLUTIONS」、arXiv:1409.4842、2014

【 0 0 1 6 】

文書8 - S.Ioffe, およびC.Szegedy, 「BATCH NORMALIZATION: ACCELERATING DEEP NETWORK TRAINING BY REDUCING INTERNAL COVARIATE SHIFT」、arXiv:1502.03167、2015

【 0 0 1 7 】

文書9 - J.M.Wolterink, T.Leiner, M.A.Viergever, およびI.Isgum, 「DILATED CONVOLUTIONAL NEURAL NETWORKS FOR CARDIOVASCULAR MR SEGMENTATION IN CONGENITAL HEART DISEASE」、arXiv:1704.03669、2017

【 0 0 1 8 】

文書10 - L.C.Piqueras, 「AUTOREGRESSIVE MODEL BASED ON A DEEP CONVOLUTIONAL NEURAL NETWORK FOR AUDIO GENERATION」、Tampere University of Technology、2016

【 0 0 1 9 】

文書11 - J.Wu, 「Introduction to Convolutional Neural Networks」、Nanjing University、2017

【 0 0 2 0 】

文書12 - I.J.Goodfellow, D.Warde-Farley, M.Mirza, A.Courville, およびY.Bengio, 「CONVOLUTIONAL NETWORKS」、Deep Learning、MIT Press、2016

【 0 0 2 1 】

文書13 - J.Gu, Z.Wang, J.Kuen, L.Ma, A.Shahrourdy, B.Shuai, T.Liu, X.Wang, およびG.Wang, 「RECENT ADVANCES IN CONVOLUTIONAL NEURAL NETWORKS」、arXiv:1512.07108、2017

【 0 0 2 2 】

文書1は、入力シーケンスを受け入れて入力シーケンス中のエントリをスコアリングする出力シーケンスを生成するために、同じ畳み込みウィンドウサイズを有する畳み込みフィルタ、パッチ正規化層、正規化線形ユニット(ReLUと省略される)層、次元変換層、指数関数的に増大する膨張畳み込み率(atrous convolution rate)を伴う膨張畳み込み層、スキップ接続、およびソフトマックス分類層を伴う、残差ブロックのグループを使用する深層畳み込みニューラルネットワークアーキテクチャを説明する。開示される技術は、文書

10

20

30

40

50

1において説明されるニューラルネットワークコンポーネントおよびパラメータを使用する。一実装形態では、開示される技術は、文書1において説明されるニューラルネットワークコンポーネントのパラメータを修正する。たとえば、文書1とは異なり、開示される技術における膨張畳み込み率は、より低い残差ブロックグループからより高い残差ブロックグループへと非指数関数的に高まる。別の例では、文書1とは異なり、開示される技術における畳み込みウィンドウサイズは、残差ブロックのグループ間で変動する。

【0023】

文書2は、文書1において説明される深層畳み込みニューラルネットワークアーキテクチャの詳細を説明する。

【0024】

文書3は、開示される技術によって使用される膨張畳み込みを説明する。本明細書では、膨張畳み込みは「拡張畳み込み(dilated convolution)」とも呼ばれる。膨張/拡張畳み込みは、少数の訓練可能なパラメータで大きな受容野を可能にする。膨張/拡張畳み込みは、膨張畳み込み率または拡張係数とも呼ばれるあるステップを用いて入力値をスキップすることによって、カーネルがその長さより長いエリアにわたって適用されるような畳み込みである。膨張/拡張畳み込みは、畳み込み演算が実行されるときに、より長い間隔の隣り合う入力エントリ(たとえば、ヌクレオチド、アミノ酸)が考慮されるように、畳み込みフィルタ/カーネルの要素間に離隔を加える。これにより、入力における長距離のコンテキスト依存性の組み込みが可能になる。膨張畳み込みは、隣接するヌクレオチドが処理されるにつれて、部分的な畳み込み計算結果を再使用のために保存する。

【0025】

文書4は、開示される技術によって使用される残差ブロックおよび残差接続を説明する。

【0026】

文書5は、開示される技術によって使用されるスキップ接続を説明する。本明細書では、スキップ接続は「ハイウェイネットワーク」とも呼ばれる。

【0027】

文書6は、開示される技術によって使用される密接続(densely connected)畳み込みネットワークアーキテクチャを説明する。

【0028】

文書7は、開示される技術によって使用される次元変換畳み込み層およびモジュールベースの処理パイプラインを説明する。次元変換畳み込みの一例は $1 \times 1$ の畳み込みである。

【0029】

文書8は、開示される技術によって使用されるバッチ正規化層を説明する。

【0030】

文書9も、開示される技術によって使用される膨張/拡張畳み込みを説明する。

【0031】

文書10は、畳み込みニューラルネットワーク、深層畳み込みニューラルネットワーク、および膨張/拡張畳み込みを伴う深層畳み込みニューラルネットワークを含む、開示される技術によって使用され得る深層ニューラルネットワークの様々なアーキテクチャを説明する。

【0032】

文書11は、サブサンプリング層(たとえば、プーリング)および全結合層を伴う畳み込みニューラルネットワークを訓練するためのアルゴリズムを含む、開示される技術によって使用され得る畳み込みニューラルネットワークの詳細を説明する。

【0033】

文書12は、開示される技術によって使用され得る様々な畳み込み演算の詳細を説明する。

【0034】

文書13は、開示される技術によって使用され得る畳み込みニューラルネットワークの様

10

20

30

40

50





e2.txtにおいて提供されることに留意されたい。

【0056】

補足テーブル3: ヒトと他の哺乳類との間で50%を超える平均ヌクレオチド保存率を伴う遺伝子だけに制約された、一般的なヒトアレル頻度で他の種において存在するミスセンスバリエーションの枯渇率。この枯渇率は、ヒトと他の種との間で同一状態であったバリエーションを使用して、稀なバリエーション(<0.1%)と比較された一般的なバリエーション(>0.1%)におけるミスセンス:同義比に基づいて計算された。このテーブルはSupplementaryTable3.txtにおいて提供されることに留意されたい。

【0057】

補足テーブル4: 一般的なヒトアレル頻度に関連する種のペアにおいて固定された置換として存在するミスセンスバリエーションの枯渇率。この枯渇率は、ヒトと関連する種のペアとの間で同一状態であったバリエーションを使用して、稀なバリエーション(<0.1%)と比較された一般的なバリエーション(>0.1%)におけるミスセンス:同義比に基づいて計算された。このテーブルはSupplementaryTable4.txtにおいて提供されることに留意されたい。

10

【0058】

補足テーブル6: SCN2A遺伝子のドメイン固有のアノテーション。ウィルコクソンの順位和のp値は、タンパク質全体と比較した特定のドメインにおけるPrimateAIスコアの相違を示す。太字で強調されたドメインはタンパク質の約7%をカバーするが、ClinVar病原性アノテーションの大半を有する。このことは、それらのドメインに対する平均PrimateAIスコアとよく符合し、PrimateAIモデルによれば上位3つの病原性ドメインである。このテーブルはSupplementaryTable6.txtにおいて提供されることに留意されたい。

20

【0059】

補足テーブル7: 予想されるミスセンス:同義比に対するアレル頻度の影響を計算する際に使用される生カウント。同義バリエーションとミスセンスバリエーションの予想されるカウントは、変異率および遺伝子変換を考慮するためにトリヌクレオチドコンテキストを使用して、イントロン領域におけるバリエーションに基づいて計算された。このテーブルはSupplementaryTables.xlsxにおいて提供されることに留意されたい。

【0060】

補足テーブル13: 3状態の二次構造および3状態の溶媒接触性予測のための深層学習モデルを訓練するために使用されるProtein DataBank(PDB)からのタンパク質名の一覧。ラベル列は、モデル訓練の訓練/妥当性確認/検定段階においてタンパク質が使用されるかどうかを示す。このテーブルはSupplementaryTable13.txtにおいて提供されることに留意されたい。

30

【0061】

補足テーブル18: タンパク質切断変異( $p < 0.05$ )のみから計算される、DDD研究において疾患との関連について名目上有意であった605個の遺伝子の一覧。このテーブルはSupplementaryTable18.txtにおいて提供されることに留意されたい。

【0062】

補足テーブル20: 少なくとも1つの観察されるDNMを伴うすべての遺伝子に対する、遺伝子ごとのde novo変異(DNM)のエンリッチメントの検定結果。すべてのDNMを含むときの、およびPrimateAIスコアが0.803より小さいミスセンスDNMを除去した後の、P値が与えられる。FDRで訂正されたP値が同様に与えられる。DDDコホートだけからの、および完全なメタ分析コホートからの、観察されるタンパク質切断(PTV)DNMとミスセンスDNMのカウントも含まれる。第1にすべてのミスセンスDNMを含むときの、および第2にPrimateAIスコアが0.803より小さいすべてのミスセンスDNMを除去した後の、観察され予測されるミスセンスDNMの同様のカウントも含まれる。このテーブルはSupplementaryTable20.txtおよびSupplementaryTable20Summary.txtにおいて提供されることに留意されたい。

40

【0063】

補足テーブル21: FDR<0.1である遺伝子におけるde novo変異のエンリッチメントを検定した結果。一度はすべてのミスセンスde novo変異についての、およびもう一度は損害を

50

引き起こすミスセンス変異だけについての、観察されるタンパク質切断(PTV)de novo変異のカウントと、他のタンパク質変換de novo変異のカウントとが含まれる。低スコアのミスセンス箇所を除外した後のP値と比較した、すべてのミスセンスサイトを含むときのP値が与えられる。このテーブルはSupplementaryTable21.txtにおいて提供されることに留意されたい。

【0064】

DataFileS1:他の種において存在するすべてのバリエーションの一覧。「ClinVar有意性」列は、利用可能な矛盾しないClinVarアノテーションを含む。このテーブルはDataFileS1.txtにおいて提供されることに留意されたい。

【0065】

DataFileS2:関連する種のペアからのすべての固定された置換の一覧。このテーブルはDataFileS2.txtにおいて提供されることに留意されたい。

【0066】

DataFileS3:霊長類とIBSである保留された(withheld)良性検定バリエーションの一覧。良性検定バリエーションは、1つ以上の霊長類の種とIBSである一般的ではないヒトバリエーションである。このテーブルはDataFileS3.txtにおいて提供されることに留意されたい。

【0067】

DataFileS4:保留された良性検定バリエーションと一致する、霊長類とIBSであるラベリングされていないバリエーションの一覧。ラベリングされていないバリエーションは、変異率、カバレッジの偏り、および霊長類の種とのアラインメント可能性について、良性検定バリエーションと一致する。このテーブルはDataFileS4.txtにおいて提供されることに留意されたい。

【0068】

Pathogenicity\_prediction\_model:一実装形態に従って開示される技術を可能にするPythonプログラミング言語のコード。このコードファイルはPathogenicity\_prediction\_model.txtにおいて提供されることに留意されたい。

【0069】

開示される技術の分野

開示される技術は、人工知能タイプコンピュータならびにデジタルデータ処理システムならびに知性のエミュレーションのための対応するデータ処理方法および製品(すなわち、知識ベースシステム、推論システム、知識取得システム)に関し、不確実性を伴う推論のためのシステム(たとえば、ファジー論理システム)、適応システム、機械学習システム、および人工ニューラルネットワークを含む。具体的には、開示される技術は、深層畳み込みニューラルネットワークを訓練するために深層学習ベースの技法を使用することに関する。

【背景技術】

【0070】

このセクションにおいて論じられる主題は、このセクションにおける言及の結果として、単なる従来技術であると見なされるべきではない。同様に、このセクションにおいて言及される問題、または背景として提供される主題と関連付けられる問題は、従来技術においてこれまで認識されていたと見なされるべきではない。このセクションの主題は異なる手法を表すにすぎず、それらの異なる手法自体も、特許請求される技術の実装形態に対応し得る。

【0071】

[機械学習]

機械学習では、出力変数を予測するために入力変数を使用される。入力変数はしばしば特徴量と呼ばれ、 $X=(X_1, X_2, \dots, X_k)$ と表記され、 $i=1, \dots, k$ である各 $X_i$ が特徴量である。出力変数はしばしば応答または依存変数と呼ばれ、変数 $Y_i$ により表記される。 $Y$ と対応する $X$ との関係は、次の一般的な形式で書くことができる。

$Y=f(x)+$

【0072】

10

20

30

40

50

上式において、 $f$ は特徴量 $(X_1, X_2, \dots, X_k)$ の関数であり、 $\epsilon$ はランダムな誤差の項である。この誤差の項は、 $X$ とは無関係であり、平均値が0である。

【0073】

実際には、特徴量 $X$ は、 $Y$ がなくても、または $X$ と $Y$ との厳密な関係を知らなくても入手可能である。誤差の項は平均値が0であるので、目標は $f$ を推定することである。

【0074】

【数1】

$$\hat{Y} = \hat{f}(X)$$

10

【0075】

上式において、

【0076】

【数2】

$$\hat{f}$$

【0077】

は  $\hat{f}$  の推定値であり、これはしばしばブラックボックスと見なされ、

【0078】

【数3】

20

$$\hat{f}$$

【0079】

の入力と出力の関係のみが知られていることを意味するが、なぜこれで機能するのかという疑問は答えられていないままである。

【0080】

関数

【0081】

【数4】

30

$$\hat{f}$$

【0082】

は学習を使用して発見される。教師あり学習および教師なし学習は、このタスクのための機械学習において使用される2つの方式である。教師あり学習では、ラベリングされたデータが訓練のために使用される。入力および対応する出力(=ラベル)を示すことによって

、関数

【0083】

【数5】

40

$$\hat{f}$$

【0084】

は、出力を近似するように最適化される。教師なし学習では、目標はラベリングされていないデータから隠された構造を見つけることである。このアルゴリズムは、入力データについての正確さの尺度を持たず、これにより教師あり学習と区別される。

【0085】

[ニューラルネットワーク]

50

図1Aは、複数の層を伴う全結合ニューラルネットワークの一実装形態を示す。ニューラルネットワークは、互いとの間でメッセージを交換する相互接続された人工ニューロン(たとえば、 $a_1$ 、 $a_2$ 、 $a_3$ )のシステムである。示されるニューラルネットワークは3つの入力を有し、2つのニューロンが隠れ層にあり、2つのニューロンが出力層にある。隠れ層は活性化関数 $f(\cdot)$ を有し、出力層は活性化関数 $g(\cdot)$ を有する。これらの接続は、適切に訓練されたネットワークが認識すべき画像を与られると正しく応答するように、訓練プロセスの間に調整された数値的な重み(たとえば、 $w_{11}$ 、 $w_{21}$ 、 $w_{12}$ 、 $w_{31}$ 、 $w_{22}$ 、 $w_{32}$ 、 $v_{11}$ 、 $v_{22}$ )を有する。入力層は生の入力を処理し、隠れ層は入力層と隠れ層との間の接続の重みに基づいて入力層から出力を処理する。出力層は、隠れ層から出力を取り込み、隠れ層と出力層との間の接続の重みに基づいてそれを処理する。ネットワークは、特徴検出ニューロンの複数の層を含む。各層は、前の層からの入力の異なる組合せに対応する多数のニューロンを有する。これらの層は、第1の層が入力画像データにおける基本的なパターンのセットを検出し、第2の層がパターンのパターンを検出し、第3の層がそれらのパターンのパターンを検出するように、構築される。

【0086】

遺伝学における深層学習の応用の概観は、以下の出版物において見出され得る。

- ・ T.Ching他、Opportunities And Obstacles For Deep Learning In Biology And Medicine、www.biorxiv.org:142760、2017
- ・ Angermueller C、Parnamaa T、Parts L、Stegle O、Deep Learning For Computational Biology. Mol Syst Biol. 2016;12:878
- ・ Park Y、Kellis M、2015 Deep Learning For Regulatory Genomics. Nat. Biotechnol. 33、825-826、(doi:10.1038/nbt.3313)
- ・ Min S、Lee B、およびYoon S、Deep Learning In Bioinformatics. Brief. Bioinform. bbw068 (2016)
- ・ Leung MK、DeLong A、Alipanahi B他、Machine Learning In Genomic Medicine: A Review of Computational Problems and Data Sets、2016
- ・ Libbrecht MW、Noble WS、Machine Learning Applications In Genetics and Genomics. Nature Reviews Genetics 2015;16(6):321-32

【先行技術文献】

【特許文献】

【0087】

【特許文献1】PCT特許出願第PCT/US2018/(未定)号(代理人整理番号ILLM 1000-9/IP-1612-PCT)

【特許文献2】国際特許出願第PCT/US2018/(未定)号(代理人整理番号第ILLM 1000-10/IP-1613-PCT)

【特許文献3】米国特許出願第(未定)号(代理人整理番号第ILLM 1000-5/IP-1611-US)

【特許文献4】米国特許出願第(未定)号(代理人整理番号第ILLM 1000-6/IP-1612-US)

【特許文献5】米国特許出願第(未定)号(代理人整理番号第ILLM 1000-7/IP-1613-US)

【特許文献6】国際特許出願公開第WO07010252号

【特許文献7】国際特許出願第PCTGB2007/003798号

【特許文献8】米国特許出願公開第2009/0088327号

【特許文献9】米国特許出願公開第2016/0085910号

【特許文献10】米国特許出願公開第2013/0296175号

【特許文献11】国際特許出願公開第WO 04/018497号

【特許文献12】米国特許第7057026号

【特許文献13】国際特許出願公開第WO 91/06678号

【特許文献14】国際特許出願公開第WO 07/123744号

【特許文献15】米国特許第7329492号

【特許文献16】米国特許第7211414号

【特許文献17】米国特許第7315019号

10

20

30

40

50

- 【特許文献 1 8】米国特許第7405281号
- 【特許文献 1 9】米国特許出願公開第2008/0108082号
- 【特許文献 2 0】米国特許第5641658号
- 【特許文献 2 1】米国特許出願公開第2002/0055100号
- 【特許文献 2 2】米国特許第7115400号
- 【特許文献 2 3】米国特許出願公開第2004/0096853号
- 【特許文献 2 4】米国特許出願公開第2004/0002090号
- 【特許文献 2 5】米国特許出願公開第2007/0128624号
- 【特許文献 2 6】米国特許出願公開第2008/0009420号
- 【特許文献 2 7】米国特許出願公開第2007/0099208A1号 10
- 【特許文献 2 8】米国特許出願公開第2007/0166705A1号
- 【特許文献 2 9】米国特許出願公開第2008/0280773A1号
- 【特許文献 3 0】米国特許出願第13/018255号
- 【特許文献 3 1】国際特許出願公開第WO 2014/142831号
- 【非特許文献】
- 【0 0 8 8】
- 【非特許文献 1】A.V.D.Oord、S.Dieleman、H.Zen、K.Simonyan、O.Vinyals、A.Graves、N.Kalchbrenner、A.Senior、およびK.Kavukcuoglu、「WAVENET: A GENERATIVE MODEL FOR RAW AUDIO」、arXiv:1609.03499、2016
- 【非特許文献 2】S.O.Arik、M.Chrzanowski、A.Coates、G.Diamos、A.Gibiansky、Y.Kang、X.Li、J.Miller、A.Ng、J.Raiman、S.Sengupta、およびM.Shoeybi、「DEEP VOICE: REAL-TIME NEURAL TEXT-TO-SPEECH」、arXiv:1702.07825、2017 20
- 【非特許文献 3】F.YuおよびV.Koltun、「MULTI-SCALE CONTEXT AGGREGATION BY DILATED CONVOLUTIONS」、arXiv:1511.07122、2016
- 【非特許文献 4】K.He、X.Zhang、S.Ren、およびJ.Sun、「DEEP RESIDUAL LEARNING FOR IMAGE RECOGNITION」、arXiv:1512.03385、2015
- 【非特許文献 5】R.K.Srivastava、K.Greff、およびJ.Schmidhuber、「HIGHWAY NETWORKS」、arXiv:1505.00387、2015
- 【非特許文献 6】G.Huang、Z.Liu、L.van der Maaten、およびK.Q.Weinberger、「DENSELY CONNECTED CONVOLUTIONAL NETWORKS」、arXiv:1608.06993、2017 30
- 【非特許文献 7】C.Szegedy、W.Liu、Y.Jia、P.Sermanet、S.Reed、D.Anguelov、D.Erhan、V.Vanhoucke、およびA.Rabinovich、「GOING DEEPER WITH CONVOLUTIONS」、arXiv:1409.4842、2014
- 【非特許文献 8】S.Ioffe、およびC.Szegedy、「BATCH NORMALIZATION: ACCELERATING DEEP NETWORK TRAINING BY REDUCING INTERNAL COVARIATE SHIFT」、arXiv:1502.03167、2015
- 【非特許文献 9】J.M.Wolterink、T.Leiner、M.A.Viergever、およびI.Isgum、「DILATED CONVOLUTIONAL NEURAL NETWORKS FOR CARDIOVASCULAR MR SEGMENTATION IN CONGENITAL HEART DISEASE」、arXiv:1704.03669、2017
- 【非特許文献 1 0】L.C.Piqueras、「AUTOREGRESSIVE MODEL BASED ON A DEEP CONVOLUTIONAL NEURAL NETWORK FOR AUDIO GENERATION」、Tampere University of Technology、2016 40
- 【非特許文献 1 1】J.Wu、「Introduction to Convolutional Neural Networks」、Nanjing University、2017
- 【非特許文献 1 2】I.J.Goodfellow、D.Warde-Farley、M.Mirza、A.Courville、およびY.Bengio、「CONVOLUTIONAL NETWORKS」、Deep Learning、MIT Press、2016
- 【非特許文献 1 3】J.Gu、Z.Wang、J.Kuen、L.Ma、A.Shahroudy、B.Shuai、T.Liu、X.Wang、およびG.Wang、「RECENT ADVANCES IN CONVOLUTIONAL NEURAL NETWORKS」、arXiv:1512.07108、2017
- 【非特許文献 1 4】T.Ching他、Opportunities And Obstacles For Deep Learning In Bi 50

ology And Medicine、www.biorxiv.org:142760、2017

【非特許文献 1 5】Angermueller C、Parnamaa T、Parts L、Stegle O、Deep Learning For Computational Biology. Mol Syst Biol. 2016;12:878

【非特許文献 1 6】Park Y、Kellis M、2015 Deep Learning For Regulatory Genomics. Nat. Biotechnol. 33、825-826、(doi:10.1038/nbt.3313)

【非特許文献 1 7】Min S、Lee B、およびYoon S、Deep Learning In Bioinformatics. Brief. Bioinform. bbw068 (2016)

【非特許文献 1 8】Leung MK、DeLong A、Alipanahi B他、Machine Learning In Genomic Medicine: A Review of Computational Problems and Data Sets、2016

【非特許文献 1 9】Libbrecht MW、Noble WS、Machine Learning Applications In Genetics and Genomics. Nature Reviews Genetics 2015;16(6):321-32 10

【非特許文献 2 0】K.He、X.Zhang、S.Ren、およびJ.Sun、「DEEP RESIDUAL LEARNING FOR IMAGE RECOGNITION」、arXiv:1512.03385、2015

【非特許文献 2 1】Bentley他、Nature 456:53-59(2008)

【非特許文献 2 2】Lizardi他、Nat.Genet.19:225-232(1998)

【非特許文献 2 3】Dunn、TamsenおよびBerry、GwennおよびEmig-Agius、DorotheaおよびJiang、YuおよびIyer、AnitaおよびUdar、NitinおよびStromberg、Michael、2017、Pisces: An Accurate and Versatile Single Sample Somatic and Germline Variant Caller、595-595、10.1145/3107411.3108203

【非特許文献 2 4】T Saunders、ChristopherおよびWong、WendyおよびSwamy、SajaniおよびBecq、JenniferおよびJ Murray、LisaおよびCheetham、Keira、2012、Strelka: Accurate somatic small-variant calling from sequenced tumor-normal sample pairs、Bioinformatics (Oxford, England)、28、1811-7、10.1093/bioinformatics/bts271 20

【非特許文献 2 5】Kim、S、Scheffler、K、Halpern、A.L.、Bekritsky、M.A、Noh、E、Kallberg M、Chen、X、Beyter、D、Krusche、P、およびSaunders、C.T、2017、Strelka2: Fast and accurate variant calling for clinical sequencing applications

【非特許文献 2 6】Stromberg、MichaelおよびRoy、RajatおよびLajugie、JulienおよびJiang、YuおよびLi、HaochenおよびMargulies、Elliott、2017、Nirvana: Clinical Grade Variant Annotator、596-596、10.1145/3107411.3108204において説明される、Illumina Inc. 30

【発明の概要】

【課題を解決するための手段】

【0089】

図面において、同様の参照文字は一般に様々な図全体で同様の部分を指す。また、図面は必ずしも縮尺通りではなく、代わりに、開示される技術の原理を示す際に一般に強調が行われる。以下の説明では、開示される技術の様々な実装形態が、以下の図面を参照して説明される。

【図面の簡単な説明】

【0090】

【図 1 A】複数の層を伴うフィードフォワードニューラルネットワークの一実装形態を示す図である。 40

【図 1 B】畳み込みニューラルネットワークの動作の一実装形態の図である。

【図 1 C】開示される技術の一実装形態による畳み込みニューラルネットワークを訓練するブロック図である。

【図 1 D】開示される技術の一実装形態によるサブサンプリング層(平均/最大プーリング)の一実装形態の図である。

【図 1 E】開示される技術の一実装形態によるReLU非線形層の一実装形態を示す図である。

【図 1 F】畳み込み層の2層の畳み込みの一実装形態を示す図である。

【図 1 G】特徴マップの追加を介して以前の情報ダウンストリームを再注入する残差接続 50

を示す図である。

【図 1 H】残差ブロックおよびスキップ接続の一実装形態を示す図である。

【図 1 I】バッチ正規化フォワードパスを示す図である。

【図 1 J】検定時のバッチ正規化変換を示す図である。

【図 1 K】バッチ正規化バックワードパスを示す図である。

【図 1 L】畳み込み層または密結合層の後と前のバッチ正規化層の使用を示す図である。

【図 1 M】1D畳み込みの一実装形態を示す図である。

【図 1 N】グローバル平均プーリング(GAP)がどのように機能するかを示す図である。

【図 1 O】拡張畳み込みを示す図である。

【図 1 P】積層(stacked)拡張畳み込みの一実装形態を示す図である。

10

【図 1 Q】開示される技術を動作させることができる例示的なコンピューティング環境を示す図である。

【図 2】本明細書で「PrimateAI」と呼ばれる、病原性予測のための深層残差ネットワークの例示的なアーキテクチャを示す図である。

【図 3】病原性分類のための深層学習ネットワークアーキテクチャであるPrimateAIの概略図である。

【図 4 A】病原性予測深層学習モデルPrimateAIの例示的なモデルアーキテクチャの詳細を示す補足テーブル16である。

【図 4 B】病原性予測深層学習モデルPrimateAIの例示的なモデルアーキテクチャの詳細を示す補足テーブル16である。

20

【図 4 C】病原性予測深層学習モデルPrimateAIの例示的なモデルアーキテクチャの詳細を示す補足テーブル16である。

【図 5】タンパク質の二次構造および溶媒接触性を予測するために使用される深層学習ネットワークアーキテクチャを示す図である。

【図 6】タンパク質の二次構造および溶媒接触性を予測するために使用される深層学習ネットワークアーキテクチャを示す図である。

【図 7 A】3状態二次構造予測深層学習(DL)モデルの例示的なモデルアーキテクチャの詳細を示す補足テーブル11である。

【図 7 B】3状態二次構造予測深層学習(DL)モデルの例示的なモデルアーキテクチャの詳細を示す補足テーブル11である。

30

【図 8 A】3状態溶媒接触性予測深層学習モデルの例示的なモデルアーキテクチャの詳細を示す補足テーブル12である。

【図 8 B】3状態溶媒接触性予測深層学習モデルの例示的なモデルアーキテクチャの詳細を示す補足テーブル12である。

【図 9】良性バリエーションおよび病原性バリエーションから基準タンパク質配列および代替タンパク質配列を生成することの一実装形態を示す図である。

【図 10】基準タンパク質配列と代替タンパク質配列をアラインメントすることの一実装形態を示す図である。

【図 11】位置特異的重み行列(PWMと省略される)または位置特異的スコアリング行列(PSSMと省略される)と呼ばれる、位置特異的頻度行列(PFMと省略される)を生成することの一実装形態を示す図である。

40

【図 12】二次構造および溶媒接触性サブネットワークの処理を示す図である。

【図 13】二次構造および溶媒接触性サブネットワークの処理を示す図である。

【図 14】二次構造および溶媒接触性サブネットワークの処理を示す図である。

【図 15】二次構造および溶媒接触性サブネットワークの処理を示す図である。

【図 16】バリエーション病原性分類器の動作を示す図である。本明細書では、バリエーションという用語は、一塩基多型(SNPと省略される)も指し、一般に一塩基バリエーション(SNVと省略される)を指す。

【図 17】残差ブロックを示す図である。

【図 18】二次構造および溶媒接触性サブネットワークのニューラルネットワークアーキ

50



テクチャを示す図である。

【図 1 9】バリエーション病原性分類器のニューラルネットワークアーキテクチャを示す図である。

【図 2 0】重要な機能ドメインのためにアノテートされた、SCN2A遺伝子の中の各アミノ酸の場所における予測される病原性スコアを示す図である。

【図 2 1 D】訓練を保留された10000個の一般的な霊長類のバリエーションの検定セットに対する良性の結果を予測することにおける分類器の比較を示す図である。

【図 2 1 E】Deciphering Developmental Disorders(DDD)の患者において発生するde novoミスセンスバリエーションに対するPrimateAI予測スコアの分布を、影響を受けていない兄弟と比較して、対応するウィルコクソンの順位和のP値とともに示す図である。

【図 2 1 F】DDD症例群vs対照群におけるde novoミスセンスバリエーションを分離する際における分類器の比較を示す図である。ウィルコクソンの順位和検定のP値が各分類器に対して示されている。

【図 2 2 A】de novoタンパク質切断変異( $P < 0.05$ )に対して有意であった605個の関連する遺伝子内での、DDDコホートからの影響を受けている個人における予想を超えるde novoミスセンス変異のエンリッチメントを示す図である。

【図 2 2 B】605個の関連する遺伝子内での、DDD患者vs影響を受けていない兄弟において発生するde novoミスセンスバリエーションに対するPrimateAI予測スコアの分布を、対応するウィルコクソンの順位和のP値とともに示す図である。

【図 2 2 C】605個の遺伝子内での症例群vs対照群におけるde novoミスセンスバリエーションを分離する際の様々な分類器の比較を示す図である。

【図 2 2 D】各分類器に対して示される曲線下面積(AUC)とともに、受信者動作特性曲線上で示される、様々な分類器の比較を示す図である。

【図 2 2 E】各分類器に対する分類の正確さおよび曲線下面積(AUC)を示す図である。

【図 2 3 A】訓練のために使用されるデータの分類の正確さに対する影響を示す図である。

【図 2 3 B】訓練のために使用されるデータの分類の正確さに対する影響を示す図である。

【図 2 3 C】訓練のために使用されるデータの分類の正確さに対する影響を示す図である。

【図 2 3 D】訓練のために使用されるデータの分類の正確さに対する影響を示す図である。

【図 2 4】一般的な霊長類バリエーションの確認に対するシーケンシングカバレッジの影響を訂正することを示す図である。

【図 2 5 A】開示されるニューラルネットワークによるタンパク質モチーフの認識を示す図である。

【図 2 5 B】開示されるニューラルネットワークによるタンパク質モチーフの認識を示す図である。

【図 2 5 C】開示されるニューラルネットワークによるタンパク質モチーフの認識を示す図である。

【図 2 6】開示されるニューラルネットワークによるタンパク質モチーフの認識を示す図である。バリエーションに対する予測される深層学習スコアへの、バリエーションの中および周りの各場所を摂動させることの影響を示す線プロットを含む。

【図 2 7】重みの相関パターンがBLOSUM62スコア行列およびGranthamスコア行列に倣っていることを示す図である。

【図 2 8 A】深層学習ネットワークのPrimateAIおよび他の分類器の性能評価を示す図である。

【図 2 8 B】深層学習ネットワークのPrimateAIおよび他の分類器の性能評価を示す図である。

【図 2 8 C】深層学習ネットワークのPrimateAIおよび他の分類器の性能評価を示す図で

10

20

30

40

50

ある。

【図 2 9 A】4つの分類器の予測スコアの分布を示す図である。

【図 2 9 B】4つの分類器の予測スコアの分布を示す図である。

【図 3 0 A】605個の疾患関連遺伝子において病原性バリエーションと良性バリエーションとを分離する際のPrimateAIネットワークおよび他の分類器の正確さを比較する図である。

【図 3 0 B】605個の疾患関連遺伝子において病原性バリエーションと良性バリエーションとを分離する際のPrimateAIネットワークおよび他の分類器の正確さを比較する図である。

【図 3 0 C】605個の疾患関連遺伝子において病原性バリエーションと良性バリエーションとを分離する際のPrimateAIネットワークおよび他の分類器の正確さを比較する図である。

【図 3 1 A】専門家により精選されたClinVarバリエーションに対する分類器の性能と、経験的なデータセットに対する性能との相関を示す図である。

10

【図 3 1 B】専門家により精選されたClinVarバリエーションに対する分類器の性能と、経験的なデータセットに対する性能との相関を示す図である。

【図 3 2】Protein Databankからのアノテートされたサンプルに対する3状態二次構造予測モデルおよび3状態溶媒接触性予測モデルの性能を示す補足テーブル14である。

【図 3 3】DSSPデータベースからのヒトタンパク質のアノテートされた二次構造ラベルを使用した深層学習ネットワークの性能比較を示す補足テーブル15である。

【図 3 4】評価した20個の分類器の各々に対する、10000個の保留された霊長類バリエーションに対する正確さの値と、DDD症例群vs対照群におけるde novoバリエーションに対するp値とを示す、補足テーブル17である。

20

【図 3 5】605個の疾患関連遺伝子に制約された、DDD症例群データセットvs対照群データセットにおけるde novoバリエーションに対する異なる分類器の性能の比較を示す補足テーブル19である。

【図 3 6】開示される半教師あり学習器のコンピューティング環境を示す図である。

【図 3 7】開示される半教師あり学習の様々なサイクルを示す図である。

【図 3 8】開示される半教師あり学習の様々なサイクルを示す図である。

【図 3 9】開示される半教師あり学習の様々なサイクルを示す図である。

【図 4 0】開示される半教師あり学習の様々なサイクルを示す図である。

【図 4 1】開示される半教師あり学習の様々なサイクルを示す図である。

【図 4 2】反復的な均衡のとれたサンプリングプロセスを示す図である。

30

【図 4 3】良性データセットを生成するために使用されるコンピューティング環境の一実装形態を示す図である。

【図 4 4】良性ヒトミスセンスSNPを生成することの一実装形態を示す図である。

【図 4 5】ヒトオーソログミスセンスSNPの一実装形態を示す図である。ヒトと一致する基準コドンおよび代替コドンを有する、ヒト以外の種におけるミスセンスSNP。

【図 4 6】ヒトと一致する基準コドンを伴うヒト以外の霊長類の種(たとえば、チンパンジー)のSNPを良性として分類することの一実装形態を示す図である。

【図 4 7】エンリッチメントスコアを計算してそれらを比較することの一実装形態を示す図である。

【図 4 8】良性SNPデータセットの一実装形態を示す図である。

40

【図 4 9 A】ヒトアレル頻度スペクトラムにわたるミスセンス:同義比を示す図である。

【図 4 9 B】ヒトアレル頻度スペクトラムにわたるミスセンス:同義比を示す図である。

【図 4 9 C】ヒトアレル頻度スペクトラムにわたるミスセンス:同義比を示す図である。

【図 4 9 D】ヒトアレル頻度スペクトラムにわたるミスセンス:同義比を示す図である。

【図 4 9 E】ヒトアレル頻度スペクトラムにわたるミスセンス:同義比を示す図である。

【図 5 0 A】他の種と同一状態であるミスセンスバリエーションに対する純化選択を示す図である。

【図 5 0 B】他の種と同一状態であるミスセンスバリエーションに対する純化選択を示す図である。

【図 5 0 C】他の種と同一状態であるミスセンスバリエーションに対する純化選択を示す図で

50

ある。

【図 5 0 D】他の種と同一状態であるミスセンスバリエーションに対する純化選択を示す図である。

【図 5 1】純化選択がない場合のヒトアレル頻度スペクトラムにわたる予想されるミスセンス:同義比を示す図である。

【図 5 2 A】CpGバリエーションおよび非CpGバリエーションに対するミスセンス:同義比を示す図である。

【図 5 2 B】CpGバリエーションおよび非CpGバリエーションに対するミスセンス:同義比を示す図である。

【図 5 2 C】CpGバリエーションおよび非CpGバリエーションに対するミスセンス:同義比を示す図である。

【図 5 2 D】CpGバリエーションおよび非CpGバリエーションに対するミスセンス:同義比を示す図である。

【図 5 3】6種の霊長類と同一状態であるヒトバリエーションのミスセンス:同義比を示す図である。

【図 5 4】6種の霊長類と同一状態であるヒトバリエーションのミスセンス:同義比を示す図である。

【図 5 5】6種の霊長類と同一状態であるヒトバリエーションのミスセンス:同義比を示す図である。

【図 5 6】調査されたヒトコホートのサイズを増やすことによって発見された新しい一般的なミスセンスバリエーションの飽和を示すシミュレーションである。

【図 5 7】ゲノムにおける異なる保存プロファイルにわたるPrimateAIの正確さを示す図である。

【図 5 8】一般的なヒトバリエーションおよびヒト以外の霊長類において存在するバリエーションからのラベリングされた良性訓練データセットへの寄与を示す補足テーブル5である。

【図 5 9】予想されるミスセンス:同義比に対するアレル頻度の影響を示す補足テーブル8である。

【図 6 0】ClinVar分析を示す補足テーブル9である。

【図 6 1】一実装形態による、ClinVarにおいて見出される他の種からのミスセンスバリエーションの数を示す補足テーブル10である。

【図 6 2】知的障害における14個の追加の遺伝子候補の発見の一実装形態を示すテーブル1である。

【図 6 3】ClinVarにおける病原性バリエーションと良性バリエーションとの間のGranthamスコアの平均の差の一実装形態を示すテーブル2である。

【図 6 4】遺伝子ごとのエンリッチメント分析の一実装形態を示す図である。

【図 6 5】ゲノムワイドエンリッチメント分析の一実装形態を示す図である。

【図 6 6】開示される技術を実装するために使用され得るコンピュータシステムの簡略化されたブロック図である。

【発明を実施するための形態】

【0091】

以下の議論は、あらゆる当業者が開示される技術を作成して使用することを可能にするために提示され、特定の適用例およびその要件の文脈で与えられる。開示される実装形態への様々な修正が当業者に容易に明らかとなり、本明細書で定義される一般的な原理は、開示される技術の趣旨および範囲から逸脱することなく他の実装形態および適用例に適用され得る。したがって、開示される技術は、示される実装形態に限定されることは意図されず、本明細書で開示される原理および特徴と矛盾しない最も広い範囲を認められるべきである。

【0092】

[ 導入 ]

[ 畳み込みニューラルネットワーク ]

10

20

30

40

50

畳み込みニューラルネットワークは特別なタイプのニューラルネットワークである。密結合層と畳み込み層との間の基本的な違いは、密層が入力特徴空間におけるグローバルパターンを学習するのに対して、畳み込み層がローカルパターンを学習するということである。画像の場合、入力の小さい2Dウィンドウにおいてパターンが見出される。この重要な特徴は、(1)畳み込みニューラルネットワークの学習するパターンが移動不変である、および(2)畳み込みニューラルネットワークがパターンの空間的階層を学習できるという、2つの興味深い特性を畳み込みニューラルネットワークに与える。

#### 【0093】

第1の特性に関して、写真の右下の角のあるパターンを学習した後、畳み込み層はそれをどこでも、たとえば左上の角において認識することができる。密結合ネットワークは、パターンが新しい位置において現れた場合、改めてパターンを学習しなければならない。これにより、畳み込みニューラルネットワークはデータ効率が高くなり、それは、一般化能力を有する表現を学習するのにより少数の訓練サンプルしか必要としないからである。

10

#### 【0094】

第2の特性に関して、第1の畳み込み層は端などの小さいローカルパターンを学習することができ、第2の畳み込み層は第1の層の特徴から作られるより大きいパターンを学習し、以下同様である。これにより、畳み込みニューラルネットワークは、ますます複雑になり抽象的になる視覚的な概念を効率的に学習することが可能になる。

#### 【0095】

畳み込みニューラルネットワークは、多くの異なる層において配置される人工ニューロンの層を、それらの層を互いに依存関係にする活性化関数を用いて相互接続することによって、高度に非線形なマッピングを学習する。畳み込みニューラルネットワークは、1つまたは複数のサブサンプリング層および非線形層とともに散在する、1つまたは複数の畳み込み層を含み、サブサンプリング層および非線形層の後には、通常は1つまたは複数の全結合層がある。畳み込みニューラルネットワークの各要素は、以前の層における特徴のセットから入力を受け取る。畳み込みニューラルネットワークは同時に学習し、それは同じ特徴マップの中のニューロンが同一の重みを有するからである。これらの局所の共有される重みがネットワークの複雑さを下げるので、多次元入力データがネットワークに入るとき、畳み込みニューラルネットワークは、特徴の抽出および回帰または分類のプロセスにおいて、データ再構築の複雑さを避ける。

20

30

#### 【0096】

畳み込みは、2つの空間軸(高さおよび幅)ならびに深さ軸(チャンネル軸とも呼ばれる)を伴う、特徴マップと呼ばれる3Dテンソルにわたって行われる。RGB画像では、深さ軸の次元は3であり、それは画像が3つの色チャンネル、すなわち赤、緑、および青を有するからである。白黒の写真では、深さは1(グレーのレベル)である。畳み込み演算は、入力特徴マップからパッチを抽出し、これらのパッチのすべてに同じ変換を適用し、出力特徴マップを生成する。この出力特徴マップはそれでも3Dテンソルであり、幅および高さを有する。その深さは任意であってよく、それは出力深さが層のパラメータであり、その深さ軸における異なるチャンネルはRGB入力におけるような特定の色をもはや表さず、むしろフィルタを表すからである。フィルタは入力データの特定の態様を符号化し、高いレベルで、単一のフィルタが、たとえば「入力における顔の存在」という概念を符号化することができる。

40

#### 【0097】

たとえば、第1の畳み込み層は、サイズ(28,28,1)の特徴マップを取り込み、サイズ(26,26,32)の特徴マップを出力する。すなわち、第1の畳み込み層は、その入力にわたる32個のフィルタを計算する。これらの32個の出力チャンネルの各々が26×26の値の格子を含み、この格子は入力にわたるフィルタの応答マップであり、入力の中の異なる位置におけるそのフィルタパターンの応答を示す。これが、特徴マップという用語が意味することである。すなわち、深さ軸におけるそれぞれの次元が特徴(またはフィルタ)であり、2Dテンソル出力[:, :, n]が入力にわたるこのフィルタの応答の2D空間マップである。

50

## 【0098】

畳み込みは、(1)通常は $1 \times 1$ 、 $3 \times 3$ 、または $5 \times 5$ である入力から抽出されたパッチのサイズ、および(2)出力特徴マップの深さという、2つの重要なパラメータによって定義され、フィルタの数は畳み込みによって計算される。しばしば、これらは32という深さで開始し、64という深さまで続き、128または256という深さで終わる。

## 【0099】

畳み込みは、3D入力特徴マップにわたってサイズ $3 \times 3$ または $5 \times 5$ のこれらのウィンドウをスライドし、それぞれの位置において止まり、周囲の特徴の3Dパッチ(形状(window\_height, window\_width, input\_depth))を抽出することによって機能する。各々のそのような3Dパッチは次いで、形状の1Dベクトル(output\_depth)への(畳み込みカーネルと呼ばれる、同じ学習された重み行列を伴うテンソル積を介して)変換される。これらのベクトルのすべてが次いで、形状の3D出力マップ(高さ、幅、output\_depth)へと空間的に再び組み立てられる。出力特徴マップの中のそれぞれの空間的位置が入力特徴マップの中の同じ位置に対応する(たとえば、出力の右下の角は入力の右下の角についての情報を含む)。たとえば、 $3 \times 3$ のウィンドウでは、ベクトル出力 $[i, j, :]$ は3Dパッチ入力 $[i-1:i+1, j-1:j+1, :]$ から来る。完全なプロセスは図1Bにおいて詳述される。

## 【0100】

畳み込みニューラルネットワークは、訓練の間に多数の勾配更新反復を介して学習される入力値と畳み込みフィルタ(重みの行列)との間で畳み込み演算を実行する、畳み込み層を備える。 $(m, n)$ をフィルタサイズとし、 $W$ は重みの行列とすると、畳み込み層は、ドット積 $w \cdot x + b$ を計算することによって、入力 $X$ を用いて $W$ の畳み込みを実行し、 $x$ は $X$ のインスタンスであり、 $b$ はバイアスである。畳み込みフィルタが入力にわたってスライドするステップサイズはストライドと呼ばれ、フィルタ面積 $(m \times n)$ は受容野と呼ばれる。同じ畳み込みフィルタが入力の異なる場所にわたって適用され、このことは学習される重みの数を減らす。このことは、すなわち、重要なパターンが入力において存在する場合、位置不変学習も可能にし、畳み込みフィルタは、重要なパターンがシーケンスの中でどこにあるかにかかわらず、重要なパターンを学習する。

## 【0101】

[畳み込みニューラルネットワークの訓練]

図1Cは、開示される技術の一実装形態による畳み込みニューラルネットワークを訓練することのブロック図を示す。畳み込みニューラルネットワークは、入力データが特定の出力推定につながるように、調整または訓練される。畳み込みニューラルネットワークは、出力推定とグラウンドトゥルースの比較に基づいて、出力推定がグラウンドトゥルースに漸近的に一致または接近するまで、逆伝播を使用して調整される。

## 【0102】

畳み込みニューラルネットワークは、グラウンドトゥルースと実際の出力との間の差に基づいてニューロン間の重みを調整することによって訓練される。これは次のように数学的に表される。

## 【0103】

【数6】

$$\Delta w_i = x_i \delta$$

## 【0104】

ただし、 $\delta = (\text{グラウンドトゥルース}) - (\text{実際の出力})$

## 【0105】

一実装形態では、訓練規則は次のように定義される。

$$w_{nm} \leftarrow w_{nm} + (t_m - o_m) \cdot o_n$$

## 【0106】

上式において、矢印は値の更新を示し、 $t_m$ はニューロン $m$ の目標値であり、 $o_m$ はニュー

10

20

30

40

50

ロン $m$ の計算された現在の出力であり、 $a_n$ は入力 $n$ であり、 $\eta$ は学習率である。

【0107】

訓練における中間ステップは、畳み込み層を使用して入力データから特徴ベクトルを生成することを含む。出力において開始して、各層における重みに関する勾配が計算される。これは、バックワードパス、または後ろに行くと呼ばれる。ネットワークにおける重みは、負の勾配および以前の重みの組合せを使用して更新される。

【0108】

一実装形態では、畳み込みニューラルネットワークは、勾配降下法によって誤差の逆伝播を実行する確率的勾配更新アルゴリズム(ADAMなど)を使用する。シグモイド関数ベースの逆伝播アルゴリズムの一例は以下のように記述される。

【0109】

【数7】

$$\varphi = f(h) = \frac{1}{1 + e^{-h}}$$

【0110】

上のシグモイド関数において、 $h$ はニューロンによって計算される加重和である。シグモイド関数は以下の導関数を有する。

【0111】

【数8】

$$\frac{\partial \varphi}{\partial h} = \varphi(1 - \varphi)$$

【0112】

このアルゴリズムは、ネットワークの中のすべてのニューロンの活性化を計算し、フォワードパスに対する出力を生み出すことを含む。隠れ層の中のニューロン $m$ の活性化は次のように記述される。

【0113】

【数9】

$$\varphi_m = \frac{1}{1 + e^{-h_m}}$$

$$h_m = \sum_{n=1}^N a_n w_{nm}$$

【0114】

これは、次のように記述される活性化を得るためにすべての隠れ層に対して行われる。

【0115】

【数10】

$$\varphi_k = \frac{1}{1 + e^{-h_k}}$$

$$h_k = \sum_{m=1}^M \varphi_m v_{mk}$$

10

20

30

40

50

【 0 1 1 6 】

そして、誤差および訂正重みが層ごとに計算される。出力における誤差は次のように計算される。

$$\delta_{ok} = (t_k - o_k) \cdot o_k(1 - o_k)$$

【 0 1 1 7 】

隠れ層における誤差は次のように計算される。

【 0 1 1 8 】

【 数 1 1 】

$$\delta_{hm} = \varphi_m(1 - \varphi_m) \sum_{k=1}^K v_{mk} \delta_{ok}$$

10

【 0 1 1 9 】

出力層の重みは次のように更新される。

$$v_{mk} = v_{mk} + \delta_{ok} \cdot a_m$$

【 0 1 2 0 】

隠れ層の重みは学習率  $\eta$  を使用して次のように更新される。

$$v_{nm} = v_{nm} + \eta \delta_{hm} a_n$$

【 0 1 2 1 】

一実装形態では、畳み込みニューラルネットワークは、すべての層にわたって誤差を計算するために勾配降下最適化を使用する。そのような最適化において、入力特徴ベクトル  $x$  および予測される出力

20

【 0 1 2 2 】

【 数 1 2 】

$$\hat{y}$$

【 0 1 2 3 】

に対して、目標が  $y$  であるときに

【 0 1 2 4 】

30

【 数 1 3 】

$$\hat{y}$$

【 0 1 2 5 】

を予測することのコストのための  $l$  として損失関数が定義され、すなわち

【 0 1 2 6 】

【 数 1 4 】

$$l(\hat{y}, y)$$

40

【 0 1 2 7 】

である。予測される出力

【 0 1 2 8 】

【 数 1 5 】

$$\hat{y}$$

【 0 1 2 9 】

は、関数  $f$  を使用して入力特徴ベクトル  $x$  から変換される。関数  $f$  は、畳み込みニューラル

50

ネットワークの重みによってパラメータ化され、すなわち

【 0 1 3 0 】

【 数 1 6 】

$$\hat{y} = f_w(x)$$

【 0 1 3 1 】

である。損失関数は

【 0 1 3 2 】

【 数 1 7 】

$$l(\hat{y}, y) = l(f_w(x), y)$$

【 0 1 3 3 】

、または $Q(z, w) = l(f_w(x), y)$ と記述され、ここで $z$ は入力データと出力データのペア $(x, y)$ である。勾配降下最適化は、以下に従って重みを更新することによって実行される。

【 0 1 3 4 】

【 数 1 8 】

$$v_{t+1} = \mu v_t - \alpha \frac{1}{n} \sum_{i=1}^N \nabla_{w_t} Q(z_i, w_t)$$

【 0 1 3 5 】

$$w_{t+1} = w_t + v_{t+1}$$

【 0 1 3 6 】

上式において、 $\alpha$ は学習率である。また、損失は $n$ 個のデータペアのセットにわたる平均として計算される。この計算は、線形収束の際に学習率が十分小さくなると終了する。他の実装形態では、計算効率をもたらすために、ネステロフの加速勾配法および適応勾配法に供給される選択されたデータペアだけを使用して、勾配が計算される。

【 0 1 3 7 】

一実装形態では、畳み込みニューラルネットワークは、コスト関数を計算するために確率的勾配降下法(SGD)を使用する。SGDは、損失関数における重みに関する勾配を、以下で記述されるように、1つのランダム化されたデータペア $z_t$ だけから計算することによって近似する。

$$v_{t+1} = \mu v_t - \alpha \nabla_{w_t} Q(z_t, w_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

【 0 1 3 8 】

上式において、 $\alpha$ は学習率であり、 $\mu$ はモメンタムであり、 $t$ は更新前の現在の重み状態である。SGDの収束速度は、学習率が十分に速く低減するときと、十分に遅く低減するときの両方において、約 $O(1/t)$ である。他の実装形態では、畳み込みニューラルネットワークは、ユークリッド損失およびソフトマックス損失などの異なる損失関数を使用する。さらなる実装形態では、Adam確率的最適化器が畳み込みニューラルネットワークによって使用される。

【 0 1 3 9 】

[ 畳み込み層 ]

畳み込みニューラルネットワークの畳み込み層は、特徴抽出器として機能する。畳み込み層は、入力データを学習して階層的特徴へと分解することが可能な、適応特徴抽出器として活動する。一実装形態では、畳み込み層は、入力として2つの画像を取り込み、出力として第3の画像を生成する。そのような実装形態では、畳み込みは2次元(2D)において2つの画像に対して動作し、一方の画像が入力画像であり、「カーネル」と呼ばれる他方の

10

20

30

40

50



画像が入力画像に対してフィルタとして適用され、出力画像を生成する。したがって、長さ $n$ の入力ベクトル $f$ および長さ $m$ のカーネル $g$ に対して、 $f$ と $g$ の畳み込み $f * g$ は次のように定義される。

【 0 1 4 0 】

【 数 1 9 】

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i - j + m/2)$$

【 0 1 4 1 】

畳み込み演算は、入力画像にわたってカーネルをスライドすることを含む。カーネルの各場所に対して、カーネルと入力画像の重複する値が乗算され、結果が加算される。この積の合計が、カーネルが中心とされる入力画像の中の点における出力画像の値である。多数のカーネルから得られた異なる出力が特徴マップと呼ばれる。

10

【 0 1 4 2 】

畳み込み層が訓練されると、それらは新しい推論データに対する認識タスクを実行するために適用される。畳み込み層は訓練データから学習するので、明示的な特徴抽出を避け、訓練データから暗黙的に学習する。畳み込み層は畳み込みフィルタカーネル重みを使用し、これは訓練プロセスの一部として決定され更新される。畳み込み層は入力の異なる特徴を抽出し、これらはより高い層において組み合わせられる。畳み込みニューラルネットワークは、様々な数の畳み込み層を使用し、それらの各々が、カーネルサイズ、ストライド、パディング、特徴マップの数、および重みなどの異なる畳み込みパラメータを伴う。

20

【 0 1 4 3 】

[ サブサンプリング層 ]

図1Dは、開示される技術の一実装形態によるサブサンプリング層の一実装形態である。サブサンプリング層は、抽出された特徴または特徴マップをノイズおよび歪みに対してロバストにするために、畳み込み層によって抽出される特徴の分解能を下げる。一実装形態では、サブサンプリング層は、2つのタイプのプーリング動作、すなわち平均プーリングおよび最大プーリングを利用する。プーリング動作は、入力を重複しない2次元空間へと分割する。平均プーリングでは、領域の中の4つの値の平均が計算される。最大プーリングでは、4つの値の最大値が選択される。

30

【 0 1 4 4 】

一実装形態では、サブサンプリング層は、その出力を最大プーリングにおける入力のうちの1つだけにマッピングし、その出力を平均プーリングにおける入力の平均にマッピングすることによる、以前の層の中のニューロンのセットに対するプーリング動作を含む。最大プーリングにおいて、プーリングニューロンの出力は、

$$\phi_o = \max(\phi_1, \phi_2, \dots, \phi_N)$$

により記述されるような、入力の中に存在する最大値である。

【 0 1 4 5 】

上式において、 $N$ はニューロンセット内の要素の総数である。

40

【 0 1 4 6 】

平均プーリングにおいて、プーリングニューロンの出力は、

【 0 1 4 7 】

【 数 2 0 】

$$\phi_o = \frac{1}{N} \sum_{n=1}^N \phi_n$$

【 0 1 4 8 】

によって記述されるような、入力ニューロンセットとともに存在する入力値の平均値であ

50

る。

【 0 1 4 9 】

上式において、Nは入力ニューロンセット内の要素の総数である。

【 0 1 5 0 】

図1Dにおいて、入力は4×4のサイズである。2×2のサブサンプリングに対して、4×4の画像は2×2のサイズの4つの重複しない行列へと分割される。平均プーリングでは、4つの値の平均は全整数出力である。最大プーリングでは、2×2の行列の中の4つの値の最大値は全整数出力である。

【 0 1 5 1 】

[ 非線形層 ]

10

図1Eは、開示される技術の一実装形態による、非線形層の一実装形態を示す。非線形層は、各隠れ層上の可能性の高い特徴の明確な識別情報をシグナリングするために、異なる非線形トリガ関数を使用する。非線形層は、正規化線形ユニット(ReLU)、双曲線正接、双曲線正接の絶対値、シグモイドおよび連続トリガ(非線形)関数を含む、非線形トリガリングを実施するために様々な固有の関数を使用する。一実装形態では、ReLU活性化は、関数 $y=\max(x,0)$ を実装し、層の入力サイズおよび出力サイズを同じに保つ。ReLUを使用することの利点は、畳み込みニューラルネットワークがより高速に多くの回数訓練されることである。ReLUは、入力が0以上の場合には、入力に関して線形であり、それ以外の場合には0である、非連続で非飽和の活性化関数である。数学的には、ReLU活性化関数は次のように記述される。

20

$$(h)=\max(h,0)$$

【 0 1 5 2 】

【 数 2 1 】

$$\varphi(h)=\begin{cases} h & \text{if } h>0 \\ 0 & \text{if } h\leq 0 \end{cases}$$

【 0 1 5 3 】

他の実装形態では、畳み込みニューラルネットワークは、

$$(h)=(a+bh)^c$$

30

によって記述される連続的な非飽和の関数である、冪ユニット活性化関数を使用する。

【 0 1 5 4 】

上式において、a、b、およびcはそれぞれ、シフト、スケール、および冪を制御するパラメータである。冪活性化関数は、cが奇数の場合にはxとyで非対称な活性化を生み出し、cが偶数の場合にはy対称な活性化を生み出すことが可能である。いくつかの実装形態では、このユニットは非正規化線形活性化を生み出す。

【 0 1 5 5 】

さらに他の実装形態では、畳み込みニューラルネットワークは、以下のロジスティック関数

【 0 1 5 6 】

40

【 数 2 2 】

$$\varphi(h)=\frac{1}{1+e^{-\beta h}}$$

【 0 1 5 7 】

によって記述される、連続的な飽和する関数である、シグモイドユニット活性化関数を使用する。

【 0 1 5 8 】

上式において、 $\beta=1$ である。シグモイドユニット活性化関数は、負の活性化を生み出さ

50

ず、y軸に関してのみ非対称である。

【 0 1 5 9 】

[ 畳み込みの例 ]

図1Fは、畳み込み層の2層の畳み込みの一実装形態を示す。図1Fにおいて、2048次元のサイズの入力が畳み込まれる。畳み込み1において、入力サイズ3×3の16個のカーネルの2つのチャンネルからなる畳み込み層によって畳み込まれる。得られる16個の特徴マップが次いで、ReLU1におけるReLU活性化関数によって正規化され、次いでサイズ3×3のカーネルを伴う16個のチャンネルプーリング層を使用して平均プーリングによってプーリング1においてプーリングされる。畳み込み2において、プーリング1の出力が次いで、3×3のサイズを伴う30個のカーネルの16個のチャンネルからなる別の畳み込み層によって畳み込まれる。さらに別のReLU2および2×2のカーネルサイズを伴うプーリング2における平均プーリングが、それに続く。畳み込み層は、可変の数、たとえば0個、1個、2個、および3個の、ストライドおよびパディングを使用する。得られる特徴ベクトルは、一実装形態によれば、512次元である。

10

【 0 1 6 0 】

他の実装形態では、畳み込みニューラルネットワークは、異なる数の畳み込み層、サブサンプリング層、非線形層、および全結合層を使用する。一実装形態では、畳み込みニューラルネットワークは、より少数の層および層当たりのより多数のニューロンを伴う浅いネットワークであり、たとえば、層当たり100個から200個のニューロンを伴う、1個、2個、または3個の全結合層である。別の実装形態では、畳み込みニューラルネットワークは、より多数の層および層当たりのより少数のニューロンを伴う深層ネットワークであり、たとえば、層当たり30個から50個のニューロンを伴う、5個、6個、または8個の全結合層である。

20

【 0 1 6 1 】

[ フォワードパス ]

特徴マップの中のf個の畳み込みコアに対するl番目の畳み込み層およびk番目の特徴マップにおける行x、列yのニューロンの出力は、次の式によって決定される。

【 0 1 6 2 】

【 数 2 3 】

$$O_{x,y}^{(l,k)} = \tanh\left(\sum_{t=0}^{f-1} \sum_{r=0}^{k_h} \sum_{c=0}^{k_w} W_{(r,c)}^{(k,t)} O_{(x+r,x+c)}^{(l-1,t)} + Bias^{(l,k)}\right)$$

30

【 0 1 6 3 】

l番目のサブサンプル層およびk番目の特徴マップにおける行x、列yのニューロンの出力は、次の式によって決定される。

【 0 1 6 4 】

【 数 2 4 】

$$O_{x,y}^{(l,k)} = \tanh\left(W^{(k)} \sum_{r=0}^{S_h} \sum_{c=0}^{S_w} O_{(x \times S_h + r, y \times S_w + c)}^{(l-1,k)} + Bias^{(l,k)}\right)$$

40

【 0 1 6 5 】

l番目の出力層のi番目のニューロンの出力は、次の式によって決定される。

【 0 1 6 6 】

【 数 2 5 】

$$O_{(l,i)} = \tanh\left(\sum_{j=0}^H O_{(l-1,j)} W_{(i,j)}^l + Bias^{(l,i)}\right)$$

50

【 0 1 6 7 】

[ 逆伝播 ]

出力層の中のk番目のニューロンの出力偏差は、次の式によって決定される。

【 0 1 6 8 】

【 数 2 6 】

$$d(O_k^o) = y_k - t_k$$

【 0 1 6 9 】

出力層の中のk番目のニューロンの入力偏差は、次の式によって決定される。

10

【 0 1 7 0 】

【 数 2 7 】

$$d(I_k^o) = (y_k - t_k)\phi'(v_k) = \phi'(v_k)d(O_k^o)$$

【 0 1 7 1 】

出力層の中のk番目のニューロンの重みおよびバイアスのばらつきは、次の式によって決定される。

【 0 1 7 2 】

【 数 2 8 】

20

$$\Delta W_{k,x}^o = d(I_k^o)y_{k,x}$$

$$\Delta Bias_k^o = d(I_k^o)$$

【 0 1 7 3 】

隠れ層の中のk番目のニューロンの出力バイアスは、次の式によって決定される。

【 0 1 7 4 】

【 数 2 9 】

30

$$d(O_k^H) = \sum_{i=0}^{i<84} d(I_i^o)W_{i,k}$$

【 0 1 7 5 】

隠れ層の中のk番目のニューロンの入力バイアスは、次の式によって決定される。

【 0 1 7 6 】

【 数 3 0 】

40

$$d(I_k^H) = \phi'(v_k)d(O_k^H)$$

【 0 1 7 7 】

隠れ層の中のk個のニューロンから入力を受け取る前の層のm番目の特徴マップの中の行x、列yにおける重みおよびバイアスのばらつきは、次の式によって決定される。

【 0 1 7 8 】

【数 3 1】

$$\Delta W_{m,x,y}^{H,k} = d(I_k^H) y_{x,y}^m$$

$$\Delta Bias_k^H = d(I_k^H)$$

【0 1 7 9】

サブサンプル層Sのm番目の特徴マップの中の行x、列yの出力バイアスは、次の式によって決定される。

10

【0 1 8 0】

【数 3 2】

$$d(O_{x,y}^{S,m}) = \sum_k^{170} d(I_{m,x,y}^H) W_{m,x,y}^{H,k}$$

【0 1 8 1】

サブサンプル層Sのm番目の特徴マップの中の行x、列yの入力バイアスは、次の式によって決定される。

【0 1 8 2】

20

【数 3 3】

$$d(I_{x,y}^{S,m}) = \phi'(v_k) d(O_{x,y}^{S,m})$$

【0 1 8 3】

サブサンプル層Sおよび畳み込み層Cのm番目の特徴マップの中の行x、列yの中の重みおよびバイアスのばらつきは、次の式によって決定される。

【0 1 8 4】

【数 3 4】

30

$$\Delta W^{S,m} = \sum_{x=0}^{fh} \sum_{y=0}^{fw} d(I_{[x/2],[y/2]}^{S,m}) O_{x,y}^{C,m}$$

$$\Delta Bias^{S,m} = \sum_{x=0}^{fh} \sum_{y=0}^{fw} d(O_{x,y}^{S,m})$$

【0 1 8 5】

畳み込み層Cのk番目の特徴マップの中の行x、列yの出力バイアスは、次の式によって決定される。

40

【0 1 8 6】

【数 3 5】

$$d(O_{x,y}^{C,k}) = d(I_{[x/2],[y/2]}^{S,k}) W^k$$

【0 1 8 7】

畳み込み層Cのk番目の特徴マップの中の行x、列yの入力バイアスは、次の式によって決定される。

【0 1 8 8】

【数 3 6】

$$d(I_{x,y}^{C,k}) = \varphi'(v_k) d(O_{x,y}^{C,k})$$

【0189】

l 番目の畳み込み層Cのk番目の特徴マップのm番目の畳み込みコアの中の行r、列cにおける重みおよびバイアスのばらつき：

【0190】

【数 3 7】

$$\Delta W_{r,c}^{k,m} = \sum_{x=0}^{f_h} \sum_{y=0}^{f_w} d(I_{x,y}^{C,k}) O_{x+r,y+c}^{l-1,m}$$

$$\Delta Bias^{C,k} = \sum_{x=0}^{f_h} \sum_{y=0}^{f_w} d(I_{x,y}^{C,k})$$

10

【0191】

[ 残差接続 ]

図1Gは、特徴マップ追加を介して以前の情報ダウンストリームを再注入する残差接続を図示する。残差接続は、過去の出力テンソルをより後の出力テンソルに追加することによって、以前の表現をデータのダウンストリームフローへと再注入することを備え、このことは、データ処理フローに沿った情報の喪失を防ぐのを助ける。残差接続は、あらゆる大規模な深層学習モデルを悩ませる2つの一般的な問題、すなわち、勾配消失および表現上のボトルネック (representational bottleneck) に対処する。一般に、10層を超える層を有するあらゆるモデルに残差接続を追加することが有益である可能性が高い。上で論じられたように、残差接続は、より前の層の出力をより後の層への入力として利用可能にして、逐次ネットワークにおけるショートカットを実質的に作成することを備える。より前の出力は、より後の活性化に連結されるのではなく、より後の活性化と加算され、このことは両方の活性化が同じサイズであると想定している。それらが異なるサイズである場合、より前の活性化を目標の形状へと再成形するための線形変換が使用され得る。残差接続についての追加の情報は、本明細書に完全に記載されるかのようにすべての目的で参照によって本明細書において引用される、K.He、X.Zhang、S.Ren、およびJ.Sun、「DEEP RESIDUAL LEARNING FOR IMAGE RECOGNITION」、arXiv:1512.03385、2015において見出され得る。

20

30

【0192】

[ 残差学習およびスキップ接続 ]

図1Hは、残差ブロックおよびスキップ接続の一実装形態を示す。残差学習の主な考え方は、残差マッピングが元のマッピングよりはるかに簡単に学習されるということである。残差ネットワークは、訓練の正確さの劣化を軽減するために、いくつかの残差ユニットを積層する。残差ブロックは、深層ニューラルネットワークにおける勾配消失をなくすために、特別な追加のスキップ接続を利用する。残差ブロックの初めにおいて、データフローは2つのストリームへと分離され、第1のストリームがブロックの変更されない入力を搬送し、一方で第2のストリームが重みおよび非線形性を適用する。ブロックの終わりにおいて、2つのストリームは要素ごとの和を使用して統合される。そのような構築の主な利点は、勾配がより簡単にネットワークを流ることが可能になることである。残差ブロックおよびスキップ接続についての追加の情報は、A.V.D.Oord、S.Dieleman、H.Zen、K.Simonyan、O.Vinyals、A.Graves、N.Kalchbrenner、A.Senior、およびK.Kavukcuoglu、「WAVENET: A GENERATIVE MODEL FOR RAW AUDIO」、arXiv:1609.03499、2016において見出され得る。

40

50

## 【0193】

残差ネットワークから利益を得て、深層畳み込みニューラルネットワーク(CNN)を簡単に訓練することができ、画像分類および物体検出の精度改善を達成することができる。畳み込みフィードフォワードネットワークは、 $l$ 番目の層の出力を $(l+1)$ 番目の層への入力として接続し、これは、以下の層遷移、すなわち $x_l = H_l(x_{l-1})$ を生じさせる。残差ブロックは、恒等関数 $x_l = H_l(x_{l-1}) + x_{l-1}$ を用いて非線形変換をバイパスするスキップ接続を追加する。残差ブロックの利点は、勾配がより後の層からより前の層へ恒等関数を通して直接流れることができるということである。しかしながら、恒等関数および $H_l$ の出力は加算によって合成され、これはネットワークにおける情報フローを妨げることがある。

## 【0194】

10

## [ 拡張畳み込み ]

図10は拡張畳み込みを示す。膨張畳み込みとも呼ばれることのある拡張畳み込みは、字面上は「穴を伴う」を意味する。フランス語の *algorithmes à trous* が名称の由来であり、これは高速二項ウェーブレット変換を計算する。これらのタイプの畳み込み層では、フィルタの受容野に対応する入力は隣り合う点ではない。これが図10に示されている。入力間の距離は拡張係数に依存する。

## 【0195】

## [ WaveNet ]

WaveNetは、生のオーディオ波形を生成するための深層ニューラルネットワークである。WaveNetは他の畳み込みネットワークから区別され、それは、WaveNetは低コストで比較的大きな「視覚野」を取り込むことが可能であるからである。その上、信号の条件をローカルおよびグローバルに追加することが可能であり、これにより、WaveNetが複数の声を伴うテキストツースピーチ(TTS)エンジンとして使用されることが可能になり、TTSはローカル条件および特定の声およびグローバル条件を与える。

20

## 【0196】

WaveNetの主なビルディングブロックは、因果的拡張畳み込みである。因果的拡張畳み込みの延長として、WaveNetは、図1Pに示されるようなこれらの畳み込みの積層を可能にする。この図において拡張畳み込みを用いて同じ受容野を取得するには、別の拡張層が必要である。積層は拡張畳み込みの反復であり、拡張畳み込み層の出力を単一の出力に接続する。これにより、WaveNetが比較的低い計算コストで1つの出力ノードの大きな「視覚野」を得ることが可能になる。比較のために、512個の入力の視覚野を得るには、完全畳み込みネットワーク(FCN)は511個の層を必要とする。拡張畳み込みネットワークの場合、8個の層が必要である。積層された拡張畳み込みは、2層の積層では7個の層、または4個の積層では6個の層しか必要ではない。同じ視覚野をカバーするために必要な計算能力の差の考え方を得るために、以下の表は、層当たり1つのフィルタおよび2というフィルタ幅という仮定のもとで、ネットワークにおいて必要とされる重みの数を示す。さらに、ネットワークが8ビットのバイナリ符号化を使用していることが仮定される。

30

## 【0197】

## 【表1】

40

ネットワークタイプ	積層数	チャンネル当たりの重みの数	重みの総数
FCN	1	$2.6 \cdot 10^5$	$2.6 \cdot 10^6$
WN	1	1022	8176
WN	2	1022	8176
WN	4	508	4064

50

## 【0198】

WaveNetは、残差接続が行われる前にスキップ接続を追加し、これはすべての後続の残差ブロックをバイパスする。これらのスキップ接続の各々は、それらを一連の活性化関数および畳み込みに通す前に加算される。直観的には、これは各層において抽出される情報の合計である。

## 【0199】

## [ バッチ正規化 ]

バッチ正規化は、データ標準化をネットワークアーキテクチャの必須の部分にすることによって、深層ネットワーク訓練を加速するための方法である。バッチ正規化は、訓練の間に時間とともに平均および分散が変化しても、データを適応的に正規化することができる。バッチ正規化は、訓練の間に見られるデータのバッチごとの平均と分散の指数移動平均を内部的に維持することによって機能する。バッチ正規化の主な影響は、残差接続とよく似て、勾配伝播を助けるので、深層ネットワークを可能にするということである。一部の超深層ネットワークは、複数のバッチ正規化層を含む場合にのみ訓練することができる。バッチ正規化についての追加の情報は、本明細書に完全に記載されるかのようによつての目的で参照によつて本明細書において引用される、S.IoffeおよびC.Szegedy、「BATCH NORMALIZATION: ACCELERATING DEEP NETWORK TRAINING BY REDUCING INTERNAL COVARIATE SHIFT」、arXiv:1502.03167、2015において見出され得る。

10

## 【0200】

バッチ正規化は、全結合層または畳み込み層のように、モデルアーキテクチャへと挿入され得るさらに別の層として見ることができる。バッチ正規化層は通常、畳み込み層または密結合層の後で使用される。バッチ正規化層は、畳み込み層または密結合層の前でも使用され得る。両方の実装形態が、開示される技術によつて使用されることが可能であり、図1Lにおいて示されている。バッチ正規化層は軸引数を取り込み、軸引数は正規化されるべき特徴軸を指定する。この引数はデフォルトでは1であり、これは入力テンソルにおける最後の軸である。これは、data\_formatが「channels\_last」に設定された状態でDense層、Conv1D層、RNN層、およびConv2D層を使用するときの正しい値である。しかし、data\_formatが「channels\_first」に設定されるConv2D層のニッチな使用事例では、特徴軸は軸1であり、バッチ正規化における軸引数は1に設定され得る。

20

## 【0201】

バッチ正規化は、入力をフィードフォワードすることと、バックワードパスを介してパラメータに関する勾配およびそれ自体の入力(its own input)を計算することとのための定義を提供する。実際には、バッチ正規化層は、畳み込み層または全結合層の後に挿入されるが、それは出力が活性化関数へと供給される前である。畳み込み層では、異なる位置にある同じ特徴マップの異なる要素、すなわち活性化が、畳み込みの性質に従うために同じ方法で正規化される。したがって、ミニバッチの中のすべての活性化は、活性化ごとではなく、すべての位置にわたって正規化される。

30

## 【0202】

内部的な共変量のシフトは、深層アーキテクチャが訓練するのに時間がかかることで悪名高かった主な理由である。これは、深層ネットワークが各層において新しい表現を学習しなければならないだけでなく、それらの分布の変化も考慮しなければならないという事実によるものである。

40

## 【0203】

一般に共変量シフトは、深層学習の領域における既知の問題であり、現実世界の問題において頻繁に発生する。よくある共変量シフト問題は、最適ではない一般化性能につながり得る、訓練セットと検定セットでの分布の違いである。この問題は通常、標準化または白色化前処理ステップ(whitening preprocessing step)によつて対処される。しかしながら、特に白色化動作は、計算負荷が高いので、共変量シフトが様々な層全体で発生する場合には特に、オンラインの状況では非現実的である。

## 【0204】

50



内部的な共変量シフトは、ネットワーク活性化の分布が、訓練の間のネットワークパラメータの変化により複数の層にわたって変化するという現象である。理想的には、各層は、各層が同じ分布を有するが機能的な関係は同じままであるような空間へと変換されるべきである。それぞれの層およびステップにおいてデータを脱相関および白色化するための、共分散行列の高価な計算を避けるために、各ミニバッチにわたる各層における各入力特徴量の分布を、平均が0になり標準偏差が1になるように正規化する。

【0205】

フォワードパス

フォワードパスの間、ミニバッチの平均および分散が計算される。これらのミニバッチの統計により、データは、平均を差し引き、標準偏差で除算することによって正規化される。最後に、データは、学習されたスケールおよびシフトパラメータを用いて、スケールリングおよびシフトされる。バッチ正規化フォワードパス $f_{BN}$ が図11に図示されている。

10

【0206】

図11において、それぞれ、 $\mu_B$ はバッチ平均であり、

【0207】

【数38】

$$\sigma_B^2$$

【0208】

はバッチ分散である。学習されたスケールおよびシフトパラメータは、それぞれ およびと表記される。分かりやすくするために、バッチ正規化手順は、本明細書では活性化ごとに説明され、対応するインデックスを省略する。

20

【0209】

正規化は微分可能な変換であるので、誤差はこれらの学習されたパラメータへと伝播され、したがって、恒等変換を学習することによってネットワークの再現能力を復元することが可能である。逆に、対応するバッチ統計と同一のスケールおよびシフトパラメータを学習することによって、それが実行すべき最適な操作であった場合、バッチ正規化変換はネットワークに対する効果を持たない。検定時に、バッチ平均および分散はそれぞれの母集団の統計により置き換えられ、それは、入力がミニバッチからの他のサンプルに依存しないからである。別の方法は、訓練の間にバッチ統計の平均をとり続け、検定時にこれらを使用してネットワーク出力を計算することである。検定時において、バッチ正規化変換は、図1Jに示されるように表現され得る。図1Jにおいて、 $\mu_D$ および

30

【0210】

【数39】

$$\sigma_D^2$$

【0211】

は、バッチ統計ではなく、それぞれ母集団の平均および分散を示す。

40

【0212】

バックワードパス

正規化は微分可能な演算であるので、バックワードパスは図1Kに図示されるように計算され得る。

【0213】

[1D畳み込み]

1D畳み込みは、図1Mに示されるように、ローカルの1Dパッチまたはサブ配列を配列から抽出する。1D畳み込みは、入力配列の中の時間的パッチから各出力タイムステップを取得する。1D畳み込み層は、配列の中のローカルパターンを認識する。同じ入力変換がバッチごとに実行されるので、入力配列の中のある場所において学習されるパターンは、異なる

40

場所においてより後に認識されることが可能であり、このことは、1D畳み込み層変換を時間的変換に対して不変にする。たとえば、サイズ5の畳み込みウィンドウを使用して塩基の配列を処理する1D畳み込み層は、長さ5以下の塩基配列を学習することが可能であるべきであり、入力配列の中の任意の文脈において塩基のモチーフを認識することが可能であるべきである。したがって、塩基レベルの1D畳み込みは、塩基の形態について学習することが可能である。

#### 【0214】

##### [グローバル平均プーリング]

図1Nは、グローバル平均プーリング(GAP)がどのように機能するかを示す。グローバル平均プーリングは、スコアリングのために最後の層の中の特徴量の空間的な平均をとることによって、分類のための全結合(FC)層を置換するために使用され得る。これは、訓練負荷を低減し、過剰適合の問題をバイパスする。グローバル平均プーリングは、モデルの前に構造的を適用し、これはあらかじめ定められた重みを伴う線形変換と等価である。グローバル平均プーリングは、パラメータの数を減らし、全結合層をなくす。全結合層は通常、最もパラメータと接続の多い層であり、グローバル平均プーリングは、同様の結果を達成するにはるかに低コストの手法を提供する。グローバル平均プーリングの主な考え方は、スコアリングのために各々の最後の層の特徴マップからの平均値を信頼性係数として生成し、直接ソフトマックス層に供給することである。

10

#### 【0215】

グローバル平均プーリングは、(1)グローバル平均プーリング層の中に余剰のパラメータがないので、グローバル平均プーリング層において過剰適合が避けられる、(2)グローバル平均プーリングの出力は特徴マップ全体の平均であるので、グローバル平均プーリングは空間的な変換に対してよりロバストになる、および(3)ネットワーク全体のすべてのパラメータの50%超を通常は占める、全結合層の中の大量のパラメータにより、それらをグローバル平均プーリング層で置き換えることで、モデルのサイズを大きく低減することができ、これがグローバル平均プーリングをモデル圧縮において非常に有用なものにする、という3つの利点を有する。

20

#### 【0216】

最後の層の中のものより強い特徴がより高い平均値を有することが予想されるので、グローバル平均プーリングは理にかなっている。いくつかの実装形態では、グローバル平均プーリングは、分類スコアのための代理として使用され得る。グローバル平均プーリングのもとでの特徴マップは、信頼性マップとして解釈されることが可能であり、特徴マップとカテゴリとの間の対応付けを強制することができる。グローバル平均プーリングは、最後の層の特徴が直接分類のために十分抽象化されている場合、特に有効であり得る。しかしながら、グローバル平均プーリング自体は、マルチレベル特徴が部分モデルのようなグループへと合成されるべきである場合には十分ではなく、これは、グローバル平均プーリングの後に、単純な全結合層または他の分類器を追加することによって最良に実行される。

30

#### 【0217】

##### [ゲノミクスにおける深層学習]

遺伝的変異は、多くの疾患の説明を助け得る。ヒトはそれぞれが固有の遺伝コードを持ち、個人のグループ内には多くの遺伝的バリエーションがある。有害な遺伝的バリエーションの大半は、自然選択によってゲノムから枯渇している。どの遺伝的変異が病原性または有害である可能性が高いかを特定することが重要である。このことは、研究者が、病原性である可能性が高い遺伝的バリエーションに注目し、多くの疾患の診断および治療を加速させることを助けるであろう。

40

#### 【0218】

バリエーションの性質および機能的な影響(たとえば、病原性)をモデル化することは重要であるが、ゲノミクスの分野においては難しい仕事である。機能的ゲノムシーケンシング技術の急速な進化にもかかわらず、バリエーションの機能的な結果の解釈には、細胞タイプに固有の転写制御システムの複雑さが原因で、大きな困難が立ちはだかっている。

50

## 【0219】

過去数十年にわたる生化学技術の進化は、これまでよりもはるかに低いコストでゲノムデータを高速に生成する、次世代シーケンシング(NGS)プラットフォームをもたらした。そのような圧倒的に大量のシーケンシングされたDNAは、アノテーションが困難なままである。教師あり機械学習アルゴリズムは通常、大量のラベリングされたデータが利用可能であるときには性能を発揮する。バイオインフォマティクスおよび多くの他のデータリッチな訓練法では、インスタンスをラベリングするプロセスが高価である。しかしながら、ラベリングされていないインスタンスは、安価であり容易に利用可能である。ラベリングされたデータの量が比較的少なく、ラベリングされていないデータの量がかなり多いシナリオでは、半教師あり学習が、手動のラベリングに対する費用対効果の高い代替手法となる。

10

## 【0220】

バリエーションの病原性を正確に予測する深層学習ベースの病原性分類器を構築するために、半教師ありアルゴリズムを使用する機会が生じる。人間の診断バイアスがない病原性バリエーションのデータベースを得ることができる。

## 【0221】

病原性分類器に関して、深層ニューラルネットワークは、高水準の特徴を連続的にモデル化するために複数の非線形の複雑な変換層を使用する、あるタイプの人工ニューラルネットワークである。深層ニューラルネットワークは、観測される出力と予測される出力との差を搬送する逆伝播を介してフィードバックを提供し、パラメータを調整する。深層ニューラルネットワークは、大きな訓練データセット、並列および分散コンピューティングの能力、および洗練された訓練アルゴリズムが利用可能になることとともに進化してきた。深層ニューラルネットワークは、コンピュータビジョン、音声認識、および自然言語処理などの、多数の領域において大きな進化を促進してきた。

20

## 【0222】

畳み込みニューラルネットワーク(CNN)および再帰型ニューラルネットワーク(RNN)は、深層ニューラルネットワークの構成要素である。畳み込みニューラルネットワークは、畳み込み層、非線形層、およびプーリング層を備えるアーキテクチャにより、画像認識において特に成功してきた。再帰型ニューラルネットワークは、パーセプトロン、長短期メモリユニット、およびゲート付き回帰型ユニットのようなビルディングブロックの間で、巡回接続を用いて入力データの連続的情報を利用するように設計される。加えて、深層空間時間ニューラルネットワーク、多次元再帰型ニューラルネットワーク、および畳み込みオートエンコーダなどの、多くの他の新興の深層ニューラルネットワークが、限られた文脈に対して提案されている。

30

## 【0223】

深層ニューラルネットワークを訓練する目的は、各層における重みパラメータの最適化であり、このことは、最も適した階層的表現をデータから学習できるように、より単純な特徴を複雑な特徴へと徐々に合成する。最適化プロセスの単一のサイクルは次のように編成される。まず、ある訓練データセットのもとで、フォワードパスが各層の中の出力を順番に計算し、ネットワークを通じて関数信号を前に伝播させる。最後の出力層において、目的損失関数が、推論された出力と所与のラベルとの間の誤差を測定する。訓練誤差を最小にするために、バックワードパスは、連鎖律を逆伝播誤差信号に使用し、ニューラルネットワーク全体のすべての重みに関する勾配を計算する。最後に、重みパラメータは、確率的勾配降下に基づく最適化アルゴリズムを使用して更新される。一方、バッチ勾配降下は、各々の完全なデータセットに対するパラメータ更新を実行し、確率的勾配降下は、データ例の各々の小さいセットに対する更新を実行することによって確率的近似を提供する。いくつかの最適化アルゴリズムは、確率的勾配低下に由来する。たとえば、AdagradおよびAdam訓練アルゴリズムは、確率的勾配降下を実行しながら、それぞれ、各パラメータのための更新頻度および勾配のモーメントに基づいて学習率を適応的に修正する。

40

## 【0224】

50

深層ニューラルネットワークの訓練における別のコア要素は正則化であり、これは、過剰適応を避けることで良好な一般化性能を達成することを意図した戦略を指す。たとえば、重み減衰は、重みパラメータがより小さい絶対値へと収束するように、目的損失関数にペナルティ項を追加する。ドロップアウトは、訓練の間にニューラルネットワークから隠れユニットをランダムに除去し、可能性のあるサブネットワークのアンサンブルであると見なされ得る。ドロップアウトの能力を高めるために、新しい活性化関数であるmaxoutと、rnnDropと呼ばれる再帰型ニューラルネットワークのためのドロップアウトの変形が提案されている。さらに、バッチ正規化は、ミニバッチ内の各活性化のためのスカラー特徴量の正規化と、各平均および分散をパラメータとして学習することを通じた、新しい正則化方法を提供する。

10

**【0225】**

シーケンシングされたデータが多次元かつ高次元であるとする、深層ニューラルネットワークは、その広い適用可能性および高い予測能力により、バイオインフォマティクスの研究に対して高い将来性がある。畳み込みニューラルネットワークは、モチーフの発見、病原性バリエーションの特定、および遺伝子発現の推論などの、ゲノムにおける配列に基づく問題を解決するために適合されてきた。畳み込みニューラルネットワークは、DNAを研究するのに特に有用である重み共有戦略を使用し、それは、この戦略が、重大な生物学的機能を有することが推定されるDNAにおける短い反復的なローカルパターンである配列モチーフを捉えることができるからである。畳み込みニューラルネットワークの特徴は、畳み込みフィルタの使用である。精巧に設計され人間により作られた特徴に基づく従来の分類手法とは異なり、畳み込みフィルタは、生の入力データを知識の有用な表現へとマッピングする処理と類似した、特徴の適応学習を実行する。この意味で、畳み込みフィルタは一連のモチーフスキャナとして機能し、それは、そのようなフィルタのセットが、入力の中の関連するパターンを認識し、訓練手順の間にそれらを更新することが可能であるからである。再帰型ニューラルネットワークは、タンパク質またはDNA配列などの、可変の長さの連続的データにおける長距離の依存関係を捉えることができる。

20

**【0226】**

したがって、バリエーションの病原性を予測するための強力な計算モデルには、基礎科学研究と橋渡し研究の両方に対して莫大な利益があり得る。

**【0227】**

一般的な多型は、多世代の自然選択によりその健康性が試されてきた自然の実験結果を表している。ヒトのミスセンス置換と同義置換についてアレル頻度分布を比較すると、ヒト以外の霊長類の種における高いアレル頻度でのミスセンスバリエーションの存在は、そのバリエーションがヒトの集団においても自然選択を受けていることを高い信頼度で予測することを発見した。対照的に、より遠縁の種における一般的なバリエーションは、進化的な距離が長くなるにつれて、負の選択を受ける。

30

**【0228】**

配列だけを使用して臨床的なde novoミスセンス変異を正確に分類する、半教師あり深層学習ネットワークを訓練するために、ヒト以外の6種の霊長類の種からの一般的な変異を利用する。500を超える既知の種により、霊長類の系統は、有意性が知られていない大半のヒトバリエーションの影響を系統的にモデル化するのに、十分な一般的な変異を含んでいる。

40

**【0229】**

ヒト基準ゲノムには、7000万個のタンパク質を変化させる可能性のあるミスセンス置換が隠れており、それらの大半は、ヒトの健康への影響が特性把握されていない稀な変異である。これらの有意性が知られていないバリエーションは、臨床応用においてゲノム解釈の課題となっており、集団全体にわたるスクリーニングおよび個別化医療のためのシーケンシングの長期的な採用の障害である。

**【0230】**

多様なヒトの集団にわたる一般的な変異の目録を作ることが、臨床的に良性の変異を特

50

定するのに有効な戦略であるが、現代のヒトから入手可能な一般的な変異は、我々の種の遠い過去におけるボトルネック事象により限られている。ヒトとチンパンジーは99%の配列相同性を共有しており、これは、チンパンジーバリエントに対して働く自然選択が、ヒトにおいて同一状態であるバリエントの影響をモデル化することの可能性を示唆している。ヒトの集団における自然な多型に対する平均合祖時間は、種の分岐時間の一部であるので、自然に発生するチンパンジー変異は大部分が、平衡選択により維持されるハプロタイプの稀な事例を除き、ヒト変異と重複しない変異空間に及ぶ。

#### 【0231】

60706人のヒトからの集約されたエクソンデータが最近利用可能になったことで、ミスセンス変異と同義変異に対するアレル頻度スペクトラムを比較することによって、この仮説を検定することが可能になった。ExACにおけるシングルTONバリエントは、トリヌクレオチドコンテキストを使用して変異率を調整した後のde novo変異により予測される、予想される2.2:1のミスセンス:同義比とよく一致するが、より高いアレル頻度では、観察されるミスセンスバリエントの数は、自然選択による有害なバリエントの除去により減少する。アレル頻度スペクトラムにわたるミスセンス:同義比のパターンは、集団における頻度が0.1%未満であるミスセンスバリエントの大部分が軽度に有害である、すなわち、集団からの即刻の除去を保証するほど病原性が高くなく、高いアレル頻度で存在することが許容されるほど中立的でもないということを示しており、これはより限られた集団データに対する以前の観察と一致している。これらの発見は、0.1%~1%より高いアレル頻度を伴うバリエントを、平衡選択および創始者効果により引き起こされるよく記録されている少数の例外を除いて、浸透性の遺伝性疾患に対しては良性である可能性が高いものとして除去するという、診療室において広く行われている経験的な実践を支持するものである。

#### 【0232】

この分析を、一般的なチンパンジーバリエント(チンパンジー集団のシーケンシングにおいて1回よりも多く観察される)と同一状態であるヒトバリエントのサブセットについて繰り返すと、ミスセンス:同義比は、アレル頻度スペクトラムにわたって概ね一定であることを発見した。チンパンジーの集団におけるこれらのバリエントの高いアレル頻度は、これらのバリエントがチンパンジーの自然選択のふるいにすでにかかけられてきたことを示し、ヒトの集団における健康へのこれらのバリエントの中立的な影響は、ミスセンスバリエントに対する選択圧力が2つの種において高度に合致していることの注目すべき証拠を与えている。チンパンジーにおいて観察されるより低いミスセンス:同義比は、軽度に有害なバリエントの効率的な除去を可能にする先祖のチンパンジーの集団におけるより大きい実効集団サイズと一貫している。

#### 【0233】

対照的に、稀なチンパンジーバリエント(チンパンジー集団のシーケンシングにおいて1回しか観察されない)は、より高いアレル頻度において、ミスセンス:同義比のあまり大きくない低下を示す。ヒト変異データからの同一サイズのコホートをシミュレートすると、このサイズのコホートにおいて一度観察されるバリエントの64%しか、集団全体において0.1%より高いアレル頻度を有せず、それと比べて、コホートにおいて複数回見られるバリエントについては99.8%が集団全体において0.1%より高いアレル頻度を有することが推定され、これは、稀なチンパンジーバリエントのすべてが選択のふるいにかかけられたとは限らないことを示している。全体として、確認されたチンパンジーミスセンスバリエントの16%が、集団全体において0.1%未満のアレル頻度を有し、より高いアレル頻度では負の選択を受けることが推定される。

#### 【0234】

次に、他のヒト以外の霊長類の種(ボノボ、ゴリラ、オランウータン、アカゲザル、およびマーモセット)において観察される変異と同一状態であるヒトバリエントを特徴付ける。チンパンジーと同様に、少数の稀なバリエント(約5~15%)の包含によるものであると推測される高いアレル頻度におけるミスセンス変異のわずかな枯渇を除き、ミスセンス:同義比がアレル頻度スペクトラムにわたって概ね等しいことを認めた。これらの結果は、

10

20

30

40

50

ミスセンスバリエントに対する選択圧が、ヒトの祖先の系統から約3500万年前に分岐したと推定される新世界ザルまでは少なくとも、霊長類の系統内で概ね合致していることを示唆する。

#### 【0235】

他の霊長類におけるバリエントと同一状態であるヒトミスセンスバリエントは、ClinVarにおける良性の結果に対して強くエンリッチメントされる。未知のまたは矛盾するアノテーションを伴うバリエントを除いた後で、霊長類オーソログを伴うヒトバリエントは、ClinVarにおいて良性または良性の可能性が高いものとしてアノテートされる確率が約95%であり、それと比較して、ミスセンス変異全般では45%であることが観察される。ヒト以外の霊長類から病原性であるものとして分類されるClinVarバリエントの小さな割合は、健康なヒトの同様のサイズのコホートからの稀なバリエントを確認することにより観察されるであろう病原性のClinVarバリエントの割合と同程度である。大きなアレル頻度データベースの出現の前に分類を受けた、病原性であるまたは病原性である可能性が高いものとしてアノテートされたこれらのバリエントのかなりの割合が、今日では異なるように評価される可能性がある。

10

#### 【0236】

ヒトの遺伝学の分野は、ヒト変異の臨床上的影響を推論するためにモデル生物に長い間依存してきたが、大半の遺伝的に扱いやすい動物モデルまでの進化的距離が長いことで、これらの発見がヒトに対してどの程度一般化可能であるかについての懸念が生まれている。ヒトおよびより遠縁の種におけるミスセンスバリエントに対する自然選択の合致を調査するために、4種の追加の哺乳類の種(ネズミ、ブタ、ヤギ、ウシ)と2種のより遠縁の脊椎動物(ニワトリ、ゼブラフィッシュ)からの概ね一般的な変異を含めるように、霊長類の系統を超えて分析を拡張した。以前の霊長類の分析とは対照的に、進化的距離が遠い場合には特に、稀なアレル頻度と比較して一般的なアレル頻度ではミスセンス変異が顕著に枯渇していることが観察され、これは、より遠縁の種における一般的なミスセンス変異のかなりの割合が、ヒトの集団においては負の選択を受けるであろうことを示している。それでも、より遠縁の脊椎動物におけるミスセンスバリエントの観察は、良性の結果の確率を高め、それは、自然選択により枯渇した一般的なミスセンスバリエントの割合は、基準であるヒトミスセンスバリエントに対して約50%よりはるかに低い枯渇率であるからである。これらの結果と一致して、ネズミ、イヌ、ブタ、およびウシにおいて観察されたヒトミスセンスバリエントは、ClinVarにおいて良性または良性の可能性が高いものとしてアノテートされる確率が約85%であり、それと比較して、霊長類の変異に対しては95%、ClinVarデータベース全体に対しては45%であることを発見した。

20

30

#### 【0237】

様々な進化的距離にある近縁の種のペアの存在も、ヒトの集団における固定されたミスセンス置換の機能的な結果を評価するための機会を与える。哺乳類の系図上で近縁の種のペア(枝長 $<0.1$ )内で、固定されたミスセンス変異が、稀なアレル頻度と比較して一般的なアレル頻度で枯渇することが観察され、これは、複数の種にわたる固定された置換のかなりの割合が、霊長類の系統内であってもヒトにおいては非中立的であることを示している。ミスセンスの枯渇の程度の比較は、複数の種にわたる固定された置換が、同一種内の多型よりはるかに中立的ではないことを示している。興味深いことに、近縁の哺乳類間での複数の種にわたる変異は、同一種内の一般的な多型と比較して、ClinVarにおいてはさほどより病原性ではなく(良性または良性の可能性が高いものとしてアノテートされる確率が83%)、これらの変化がタンパク質の機能を無効にするのではなく、むしろ、種固有の適応的な利益を授けるタンパク質機能の調整を招いていることを示唆する。

40

#### 【0238】

有意性が知られてない多数の潜在的なバリエントがあること、および臨床上的応用には正確なバリエント分類が決定的に重要であることにより、機械学習を用いた問題の解決が多く試みられてきたが、これらの努力は、一般的なヒトバリエントの量が不十分であること、および精選されたデータベースにおけるアノテーションの品質が疑わしいことにより

50

大きく制約されてきた。6種のヒト以外の霊長類からの変異は、一般的なヒト変異と重複せず大部分が良性の結果をもたらす300000個を超える固有のミスセンスバリエーションに寄与し、機械学習手法に使用できる訓練データセットのサイズを大きく拡大した。

#### 【0239】

人間により加工された多数の特徴およびメタ分類器を利用するこれまでのモデルと異なり、対象のバリエーションの側にあるアミノ酸配列および他の種におけるオーソログな配列アラインメントのみを入力として取り込む、単純な深層学習残差ネットワークを適用する。タンパク質構造についての情報をネットワークに提供するために、配列だけから二次構造および溶媒接触性を学習するように2つの別々のネットワークを訓練し、これらをサブネットワークとしてより大きな深層学習ネットワークに組み込み、タンパク質構造に対する影響を予測する。配列を開始点として使用することで、不完全に確認されている可能性がある、または矛盾して適用されている可能性がある、タンパク質構造および機能ドメインのアノテーションにおける存在し得るバイアスが回避される。

10

#### 【0240】

良性である可能性が高い霊長類バリエーションと、変異率およびシーケンシングカバレッジについて一致するランダムな未知のバリエーションとを分離するように、ネットワークのアンサンブルを最初に訓練することによって、訓練セットが良性のラベルを持つバリエーションしか含まないという問題を克服するために、半教師あり学習を使用する。このネットワークのアンサンブルは、未知のバリエーションの完全なセットをスコアリングするために、および、より病原性であるという予測される結果を持つ未知のバリエーションに向かってバイアスをかけることによって分類器の次の反復をシードするように未知のバリエーションの選択に影響を与えるために使用され、モデルが準最適な結果へと尚早に収束するのを防ぐために各反復において緩やかなステップをとる。

20

#### 【0241】

一般的な霊長類の変異はまた、メタ分類器の増殖により客観的に評価することが難しくなっている既存の方法を評価するための、以前に使用された訓練データとは完全に無関係であるクリーンな評価データセットを提供する。10000個の提供された霊長類の一般的なバリエーションを使用して、4つの他の人気のある分類アルゴリズム(Sift、Polyphen2、CADD、M-CAP)とともに、我々のモデルの性能を評価した。すべてのヒトミスセンスバリエーションの概ね50%は、一般的なアレル頻度では自然選択によって除去されるので、変異率によって、10000個の提供された霊長類の一般的なバリエーションと一致したランダムに選ばれたミスセンスバリエーションのセットに対して、各分類器について50パーセントのスコアを計算し、その閾値を使用して、提出された霊長類の一般的なバリエーションを評価した。我々の深層学習モデルの正確さは、ヒトの一般的なバリエーションだけで訓練された深層学習ネットワークを使用しても、またはヒトの一般的なバリエーションと霊長類のバリエーションの両方を使用しても、この独立の評価データセットについて、他の分類器よりはるかに良好であった。

30

#### 【0242】

最近のトリオシーケンシング研究は、神経発達障害を持つ患者と患者の健康な兄弟における数千個のde novo変異の目録を作っており、症例群vs対照群におけるde novoミスセンス変異を分離する際の様々な分類アルゴリズムの強さの評価を可能にしている。4つの分類アルゴリズムの各々について、症例群vs対照群における各de novoミスセンスバリエーションをスコアリングし、2つの分布の間の差のウィルコクソンの順位和検定からのp値を報告し、この臨床シナリオでは、霊長類バリエーションについて訓練された深層学習方法(p約 $10^{-3}$ )が他の分類器(p約 $10^{-13}$ から $10^{-19}$ )はるかに良好な性能であったことを示した。このコホートについて以前に報告された予想を超える、de novoミスセンスバリエーションの約1.3-foldエンリッチメントから、およびミスセンスバリエーションの約20%が機能喪失の影響を生むという以前の推定から、完璧な分類器はp約 $10^{-40}$ というp値で2つのクラスを分離することが予想され、これは我々の分類器に改善の余地がまだあることを示している。

40

#### 【0243】

50

深度学习分類器の正確さは訓練データセットのサイズと符合し、6種の霊長類の各々からの変異データは独立に、分類器の正確さを上げることに寄与する。ヒト以外の霊長類の種が多数かつ多様にあることは、タンパク質を変化させるバリエーションに対する選択圧力が霊長類の系統内で概ね合致していることを示す証拠とともに、臨床上のゲノム解釈を現在制約している、有意性が知られていない数百万個のヒトバリエーションを分類するための効果的な戦略として、系統的な霊長類集団のシーケンシングを示唆する。504種の知られているヒト以外の霊長類の種のうち、約60%が狩猟および生息地喪失により絶滅に瀕しており、これらの固有の代替のいない種と我々自身の両方に利益をもたらすであろう、緊急を要する世界的な保全の努力に対する動機となっている。

【0244】

10

ゲノムデータ全体はエクソデータほど集約された形では利用可能ではないが、深いイントロン領域における自然選択の影響を検出するための能力を制限することで、エクソン領域から遠く離れた隠れたスプライシング変異の観察されるカウントと予想されるカウントを計算することも可能になった。全体として、エクソイントロン境界から50ntを超える距離にある隠れたスプライシング変異において、60%の欠失を認めた。信号の減衰は、エクソンと比較してゲノムデータ全体ではサンプルサイズがより小さいことと、深いイントロンバリエーションの影響を予測することがより難しいこととの組合せによるものである可能性が高い。

【0245】

20

[用語]

限定はされないが、特許、特許出願、論説、書籍、論文、およびウェブページを含む、本出願において引用されるすべての文献および同様の資料は、そのような文献および同様の資料のフォーマットとは無関係に、全体が参照によって明確に引用される。限定はされないが、定義される用語、用語の使用法、説明される技法などを含めて、引用される文献および同様の資料のうちの一つまたは複数が、本出願とは異なる場合、または本出願と矛盾する場合、本出願が優先する。

【0246】

本明細書では、以下の用語は示される意味を有する。

【0247】

塩基は、ヌクレオチド塩基またはヌクレオチド、すなわちA(アデニン)、C(シトシン)、T(チミン)、またはG(グアニン)を指す。

30

【0248】

本出願は、「タンパク質」および「翻訳配列」という用語を交換可能に使用する。

【0249】

本出願は、「コドン」および「塩基トリプレット」という用語を交換可能に使用する。

【0250】

本出願は、「アミノ酸」および「翻訳単位」という用語を交換可能に使用する。

【0251】

本出願は、「バリエーション病原性分類器」、「バリエーション分類のための畳み込みニューラルネットワークベースの分類器」、および「バリエーション分類のための深層畳み込みニューラルネットワークベースの分類器」という語句を交換可能に使用する。

40

【0252】

「染色体」という用語は、生きている細胞の遺伝情報を持っている遺伝子の担体を指し、これはDNAおよびタンパク質の構成要素(特にヒストン)を備えるクロマチン鎖に由来する。従来国際的に認識されている個々のヒトゲノム染色体ナンバリングシステムが本明細書で利用される。

【0253】

「サイト」という用語は、基準ゲノム上の一意な場所(たとえば、染色体ID、染色体の場所および向き)を指す。いくつかの実装形態では、サイトは、残基、配列タグ、または配列上のセグメントの場所であり得る。「座」という用語は、基準染色体上での核酸配列

50



または多型の具体的な位置を指すために使用され得る。

【0254】

本明細書の「サンプル」という用語は、典型的には、シーケンシングおよび/もしくはフェージングされるべき少なくとも1つの核酸配列を含有する核酸もしくは核酸の混合物を含有する、体液、細胞、組織、器官、または生物体に由来する、サンプルを指す。そのようなサンプルは、限定はされないが、唾液/口腔液、羊水、血液、血液の断片、細針生検サンプル(たとえば、直視下生検、細針生検など)、尿、腹膜液、胸膜液、組織外植、器官培養、および任意の他の組織もしくは細胞の標本、またはそれらの一部もしくはそれらの派生物、またはそれらから分離されたものを含む。サンプルはしばしば、ヒト対象(たとえば、患者)から取られるが、サンプルは、限定はされないが、イヌ、ネコ、ウマ、ヤギ、ヒツジ、ウシ、ブタなどを含む、染色体を有する任意の生物体から取ることができる。サンプルは、生物学的な供給源から得られるものとして直接使用されることがあり、または、サンプルの特性を修正するための前処理の後に使用されることがある。たとえば、そのような前処理は、血液から血漿を調製すること、粘液を希釈することなどを含み得る。前処理の方法はまた、限定はされないが、濾過、沈殿、希釈、蒸留、混合、遠心分離、凍結、凍結乾燥、濃縮、増幅、核酸断片化、干渉する要素の不活性化、試薬の追加、溶解などを伴い得る。

10

【0255】

「配列」という用語は、互いに結合されたヌクレオチドの鎖を含み、または表す。ヌクレオチドはDNAまたはRNAに基づき得る。1つの配列は複数の部分配列を含み得ることを理解されたい。たとえば、(たとえばPCRアンプリコン)の単一配列は350個のヌクレオチドを有し得る。サンプルリードは、これらの350個のヌクレオチド内の複数の部分配列を含み得る。たとえば、サンプルリードは、たとえば20~50個のヌクレオチドを有する、第1および第2のランキング部分配列を含み得る。第1および第2のランキング部分配列は、対応する部分配列(たとえば、40~100個のヌクレオチド)を有する反復的なセグメントの両側に位置し得る。ランキング部分配列の各々は、プライマー部分配列(たとえば、10~30個のヌクレオチド)を含み得る(またはその一部を含み得る)。読むのを簡単にするために、「部分配列」という用語は「配列」と呼ばれるが、2つの配列は必ずしも共通の鎖上で互いに別々であるとは限らないことを理解されたい。本明細書で説明される様々な配列を区別するために、配列は異なるラベル(たとえば、標的配列、プライマー配列、ランキング配列、基準配列など)を与えられ得る。「アレル」などの他の用語は、同様の物を区別するために異なるラベルを与えられ得る。

20

30

【0256】

「ペアエンドシーケンシング(paired-end sequencing)」という用語は、標的フラグメントの両端をシーケンシングするシーケンシング方法を指す。ペアエンドシーケンシングは、ゲノム再配置および反復セグメント、ならびに遺伝子融合および新規転写物の検出を容易にし得る。ペアエンドシーケンシングの方法論は、各々が本明細書において参照によって引用される、国際特許出願公開第WO07010252号、国際特許出願第PCTGB2007/003798号、および米国特許出願公開第2009/0088327号において説明されている。一例では、一連の操作は次のように実行され得る。(a)核酸のクラスタを生成する。(b)核酸を直線化する。(c)第1のシーケンシングプライマーをハイブリダイゼーションし、上で記載されたような延長、走査、およびデプロッキングの繰り返されるサイクルを実行する。(d)相補的なコピーを合成することによってフローセル表面上の標的核酸を「逆にする」。(e)再合成された鎖を直線化する。(f)第2のシーケンシングプライマーをハイブリダイゼーションし、上で記載されたような延長、走査、およびデプロッキングの繰り返されるサイクルを実行する。この逆転操作は、ブリッジ増幅の単一サイクルについて上に記載されたように試薬を導入するために実行され得る。

40

【0257】

「基準ゲノム」または「基準配列」という用語は、対象からの特定された配列の基準にするために使用され得る任意の生物体の任意の特定の既知のゲノム配列を、それが部分的

50

なものであるか完全なものであるかにかかわらず指す。たとえば、ヒト対象ならびに多くの他の生物体のために使用される基準ゲノムは、ncbi.nlm.nih.govの米国国立生物工学情報センターにおいて見つかる。「ゲノム」は、核酸配列で表現される、生物体またはウイルスの完全な遺伝情報を指す。ゲノムは、遺伝子とDNAのノンコーディング配列の両方を含む。基準配列は、それとアラインメントされるリードより大きいことがある。たとえば、それは少なくとも約100倍大きいことがあり、または少なくとも約1000倍大きいことがあり、または少なくとも約10000倍大きいことがあり、または少なくとも約105倍大きいことがあり、または少なくとも約106倍大きいことがあり、または少なくとも約107倍大きいことがある。一例では、基準ゲノム配列は、完全な長さのヒトゲノムの基準ゲノム配列である。別の例では、基準ゲノム配列は、13番染色体などの特定のヒト染色体に限定される。いくつかの実装形態では、基準染色体は、ヒトゲノムバージョンhg19からの染色体配列である。そのような配列は染色体基準配列と呼ばれ得るが、基準ゲノムという用語がそのような配列を包含することが意図される。基準配列の他の例には、他の種のゲノム、ならびに任意の種の染色体、部分染色体領域(鎖など)などがある。様々な実装形態において、基準ゲノムは、複数の個体由来するコンセンサス配列または他の組合せである。しかしながら、いくつかの適用例では、基準配列は特定の個体から取られることがある。

10

20

30

40

50

#### 【0258】

「リード」という用語は、ヌクレオチドサンプルまたは基準のフラグメントを記述する配列データの集合体を指す。「リード」という用語は、サンプルリードおよび/または基準リードを指し得る。通常、必須ではないが、リードは、サンプルまたは基準における連続的な塩基対の短い配列を表す。リードは、サンプルまたは基準フラグメントの塩基対配列によって文字で(ATCGで)表され得る。リードは、メモリデバイスに記憶され、リードが基準配列と一致するかどうか、または他の基準を満たすかどうかを決定するために適宜処理され得る。リードは、シーケンシング装置から直接、またはサンプルに関する記憶された配列情報から間接的に得られ得る。いくつかの場合、リードは、たとえば染色体またはゲノム領域または遺伝子にアラインメントされ具体的に割り当てられ得る、より大きい配列または領域を特定するために使用され得る、十分な長さの(たとえば、少なくとも約25bp)DNA配列である。

#### 【0259】

次世代シーケンシング方法には、たとえば、合成技術によるシーケンシング(Illumina)、パイロシーケンシング(454)、イオン半導体技術(Ion Torrentシーケンシング)、単一分子リアルタイムシーケンシング(Pacific Biosciences)、およびライゲーションによるシーケンシング(SOLiDシーケンシング)がある。シーケンシング方法に応じて、各リードの長さは、約30bpから10000bp以上にまで変動し得る。たとえば、SOLiDシーケンサを使用するIlluminaシーケンシング方法は、約50bpの核酸リードを生成する。別の例では、Ion Torrentシーケンシングは最高で400bpの核酸リードを生成し、454パイロシーケンシングは約700bpの核酸リードを生成し得る。さらに別の例では、単一分子リアルタイムシーケンシング方法は、10000bpから15000bpのリードを生成し得る。したがって、いくつかの実装形態では、核酸配列リードは、30~100bp、50~200bp、または50~400bpの長さを有する。

#### 【0260】

「サンプルリード」、「サンプル配列」、または「サンプルフラグメント」という用語は、サンプルからの対象のゲノム配列の配列データを指す。たとえば、サンプルリードは、フォワードプライマー配列およびリバースプライマー配列を有するPCRアンプリコンからの配列データを備える。配列データは、任意の配列選択方法から得られ得る。サンプルリードは、たとえば、sequencing-by-synthesis(SBS)反応、sequencing-by-ligation反応、または、そのために反復要素の長さおよび/または正体を決定することが望まれる任意の他の適切なシーケンシング方法からのものであり得る。サンプルリードは、複数のサンプルリードに由来するコンセンサス(たとえば、平均または加重)配列であり得る。いくつかの実装形態では、基準配列を提供することは、PCRアンプリコンのプライマー配列に基

づいて対象座を特定することを備える。

【0261】

「生フラグメント」という用語は、サンプルリードまたはサンプルフラグメント内で指定場所または二次的な対象場所と少なくとも部分的に重複する、対象のゲノム配列の部分に対する配列データを指す。生フラグメントの非限定的な例には、duplex stitchedフラグメント、simplex stitchedフラグメント、duplex un-stitchedフラグメント、およびsimplex un-stitchedフラグメントがある。「生」という用語は、生フラグメントがサンプルリードの中の潜在的なバリエーションに対応しそれが本物であることを証明または確認する、支持バリエーションを呈するかどうかにかかわらず、サンプルリードの中の配列データに対する何らかの関連を有する配列データを含むことを示すために使用される。「生フラグメント」という用語は、フラグメントが、サンプルリードの中のバリエーションコールを妥当性確認する支持バリエーションを必ず含むことを示さない。たとえば、サンプルリードが第1のバリエーションを呈することが、バリエーションコールアプリケーションによって決定されるとき、バリエーションコールアプリケーションは、1つまたは複数の生フラグメントが、サンプルリードの中にそのバリエーションがあるとすれば存在することが予想され得る対応するタイプの「支持」バリエーションを欠いていることを決定し得る。

10

【0262】

「マッピング」、「アラインメントされる」、「アラインメント」、または「アラインメントしている」という用語は、リードまたはタグを基準配列と比較し、それにより、基準配列がリード配列を含むかどうかを決定するプロセスを指す。基準配列がリードを含む場合、リードは、基準配列にマッピングされることがあり、またはいくつかの実装形態では、基準配列の中の特定の位置にマッピングされることがある。いくつかの場合、アラインメントは単に、リードが特定の基準配列のメンバーであるかどうか(すなわち、リードが基準配列の中に存在するかしないか)を伝える。たとえば、ヒト13番染色体の基準配列に対するリードのアラインメントは、リードが13番染色体の基準配列の中に存在するかどうかを伝える。この情報を提供するツールは、セットメンバーシップテスターと呼ばれ得る。いくつかの場合、アラインメントは追加で、リードまたはタグがマッピングする基準配列の中の位置を示す。たとえば、基準配列がヒトゲノム配列全体である場合、アラインメントは、リードが13番染色体上に存在することを示すことがあり、さらに、リードが13番染色体の特定の鎖および/またはサイトにあることを示すことがある。

20

30

【0263】

「インデル」という用語は、生物体のDNAにおける塩基の挿入および/または欠失を指す。マイクロインデルは、1~50個のヌクレオチドの正味の変化をもたらすインデルを表す。ゲノムのコーディング領域において、インデルの長さが3の倍数ではない限り、インデルはフレームシフト変異を生み出す。インデルは点変異と対比され得る。インデルは配列からヌクレオチドを挿入または削除するが、点変異はDNAの全体の数を変えずにヌクレオチドのうち1つを置き換えるある形式の置換である。インデルは、タンデム塩基変異(TBM)とも対比することができ、TBMは隣接するヌクレオチドにおける置換として定義され得る(主に2つの隣接するヌクレオチドにおける置換、しかし3つの隣接するヌクレオチドにおける置換が観察されている)。

40

【0264】

「バリエーション」という用語は、核酸基準と異なる核酸配列を指す。典型的な核酸配列バリエーションには、限定はされないが、一塩基多型(SNP)、短い欠失および挿入の多型(インデル)、コピー数変異(CNV)、マイクロサテライトマーカー、またはショートタンデムリピートおよび構造変異がある。体細胞バリエーションコーリング(somatic variant calling)は、DNAサンプルにおいて低頻度に存在するバリエーションを特定するための試みである。体細胞バリエーションコーリングは、癌治療の文脈において関心の対象である。癌はDNAの変異の蓄積により引き起こされる。腫瘍からのDNAサンプルは一般に異質であり、いくつかの正常細胞、癌進行の早期段階にあるいくつかの細胞(少数の変異を伴う)、およびいくつかの後段階の細胞(多数の変異を伴う)を含む。この異質さにより、(たとえば、FFPEサンプルか

50

ら)腫瘍をシーケンシングするとき、体細胞突然変異がしばしば低頻度で現れる。たとえば、ある所与の塩基を含むリードの10%だけにおいて、SNVが見られることがある。バリエーション分類器によって体細胞性または生殖細胞性であると分類されるべきバリエーションは、「検定対象バリエーション(variant under test)」とも本明細書では呼ばれる。

【0265】

「ノイズ」という用語は、シーケンシングプロセスおよび/またはバリエーションコールアプリケーションにおける1つまたは複数のエラーに起因する誤ったバリエーションコールを指す。

【0266】

「バリエーション頻度」という用語は、割合または百分率で表される、ある集団の中の特定の座におけるアレル(遺伝子のバリエーション)の相対的な頻度を表す。たとえば、この割合または百分率は、そのアレルを持つ集団の中のすべての染色体の割合であり得る。例として、サンプルバリエーション頻度は、ある個人からの対象のゲノム配列について取得されたリードおよび/またはサンプルの数に対応する「集団」にわたる、対象のゲノム配列に沿った特定の座/場所におけるアレル/バリエーションの相対的な頻度を表す。別の例として、基準バリエーション頻度は、1つまたは複数の基準ゲノム配列に沿った特定の座/場所におけるアレル/バリエーションの相対的な頻度を表し、リードおよび/またはサンプルの数に対応する「集団」は、正常な個人の集団からの1つまたは複数の基準ゲノム配列について取得される。

10

【0267】

「バリエーションアレイ頻度(VAF)」という用語は、標的場所における、バリエーションと一致することが観察されたシーケンシングされたリードをカバレッジ全体で割った百分率を指す。VAFはバリエーションを持つシーケンシングされたリードの比率の尺度である。

20

【0268】

「場所」、「指定場所」、および「座」という用語は、ヌクレオチドの配列内の1つまたは複数のヌクレオチドの位置または座標を指す。「場所」、「指定場所」、および「座」という用語は、ヌクレオチドの配列の中の1つまたは複数の塩基対の位置または座標も指す。

【0269】

「ハプロタイプ」という用語は、一緒に受け継がれる染色体上の隣接するサイトにおけるアレルの組合せを指す。ハプロタイプは、所与の座のセット間で組み換え事象が発生した場合にはその数に依存して、1つの座、いくつかの座、または染色体全体であり得る。

30

【0270】

本明細書の「閾値」という用語は、サンプル、核酸、またはその一部(たとえば、リード)を特徴付けるためにカットオフとして使用される、数値または数字ではない値を指す。閾値は経験的な分析に基づいて変動し得る。閾値は、そのような値の示唆をもたらす源がある特定の方式で分類されるべきであるかどうかを決定するために、測定された値または計算された値と比較され得る。閾値は経験的または分析的に特定され得る。閾値の選択は、ユーザが分類を行うために有することを望む信頼性のレベルに依存する。閾値は特定の目的で(たとえば、感度と選択度のバランスをとるように)選ばれ得る。本明細書では、「閾値」という用語は、分析のコースが変更され得る点、および/または活動が惹起され得る点を示す。閾値は所定の数である必要はない。代わりに、閾値は、たとえば、複数の要因に基づく関数であり得る。閾値は状況に適應するものであり得る。その上、閾値は、上限、下限、または制限値間の範囲を示し得る。

40

【0271】

いくつかの実装形態では、シーケンシングデータに基づく尺度またはスコアが閾値と比較され得る。本明細書では、「尺度」または「スコア」という用語は、シーケンシングデータから決定された値もしくは結果を含むことがあり、または、シーケンシングデータから決定された値もしくは結果に基づく関数を含むことがある。閾値と同様に、尺度またはスコアは状況に適應するものであり得る。たとえば、尺度またはスコアは正規化された値であり得る。スコアまたは尺度の例として、1つまたは複数の実装形態は、データを分析

50

するときにカウントスコアを使用し得る。カウントスコアはサンプルリードの数に基づき得る。サンプルリードは1つまたは複数のフィルタリング段階を経ていることがあるので、サンプルリードは少なくとも1つの一般的な特性または品質を有する。たとえば、カウントスコアを決定するために使用されるサンプルリードの各々は、基準配列とアラインメントされていることがあり、または潜在的なアレルとして割り当てられることがある。一般的な特性を有するサンプルリードの数はリードカウントを決定するためにカウントされ得る。カウントスコアはリードカウントに基づき得る。いくつかの実装形態では、カウントスコアはリードカウントに等しい値であり得る。他の実装形態では、カウントスコアはリードカウントおよび他の情報に基づき得る。たとえば、カウントスコアは、遺伝子座の特定のアレルに対するリードカウントおよび遺伝子座に対するリードの総数に基づき得る。いくつかの実装形態では、カウントスコアは、遺伝子座に対するリードカウントおよび以前に得られたデータに基づき得る。いくつかの実装形態では、カウントスコアは複数の所定の値の間の正規化されたスコアであり得る。カウントスコアはまた、サンプルの他の座からのリードカウントの関数、または対象サンプルと同時に実行された他のサンプルからのリードカウントの関数であり得る。たとえば、カウントスコアは、特定のアレルのリードカウントおよびサンプルの中の他の座のリードカウントおよび/または他のサンプルからのリードカウントの関数であり得る。一例として、他の座からのリードカウントおよび/または他のサンプルからのリードカウントが、特定のアレルに対するカウントスコアを正規化するために使用され得る。

10

20

30

40

50

#### 【0272】

「カバレッジ」または「フラグメントカバレッジ」という用語は、配列の同じフラグメントに対するサンプルリードの数のカウントまたは他の尺度を指す。リードカウントは対応するフラグメントをカバーするリードの数のカウントを表し得る。あるいは、カバレッジは、履歴の知識、サンプルの知識、座の知識などに基づく指定された係数を、リードカウントと乗じることによって決定され得る。

#### 【0273】

「リード深さ」(慣習的に「 $\times$ 」が後に続く数)という用語は、標的場所における重複するアラインメントを伴うシーケンシングされたリードの数を指す。これはしばしば、平均として、または間隔(エクソン、遺伝子、またはパネルなど)のセットにわたってカットオフを超える百分率として表される。たとえば、パネル平均カバレッジが $1.105\times$ であり、カバーされる標的塩基の98%が $>100\times$ であるということ、臨床報告が述べることもある。

#### 【0274】

「塩基コール品質スコア」または「Qスコア」という用語は、単一のシーケンシングされた塩基が正しい確率に反比例する、0~20の範囲のPHREDスケールされた確率を指す。たとえば、Qが20であるT塩基コールは、0.01という信頼性P値を伴い正しい可能性が高いと見なされる。Q<20であるあらゆる塩基コールは低品質であると見なされるべきであり、バリエーションを支持するシーケンシングされたリードのかなりの部分が低品質であるようなあらゆる特定されたバリエーションは、偽陽性の可能性があるから見なされるべきである。

#### 【0275】

「バリエーションリード」または「バリエーションリード数」という用語は、バリエーションの存在を支持するシーケンシングされたリードの数を指す。

#### 【0276】

#### [シーケンシングプロセス]

本明細書に記載される実装形態は、配列の変異を特定するために核酸配列を分析することに適用可能であり得る。実装形態は、遺伝子の場所/座の潜在的なバリエーション/アレルを分析し、遺伝子座の遺伝子型を決定するために、言い換えると、座に対する遺伝子型コールを提供するために使用され得る。例として、核酸配列は、米国特許出願公開第2016/0085910号および米国特許出願公開第2013/0296175号において説明される方法およびシステムに従って分析されることがあり、これらの出願公開の完全な主題の全体が、本明細書にお

いて参照によって明確に引用される。

【0277】

一実装形態では、シーケンシングプロセスは、DNAなどの核酸を含む、または含むことが疑われるサンプルを受け取ることを含む。サンプルは、動物(たとえばヒト)、植物、バクテリア、または菌類などの、既知のまたは未知の源からのものであり得る。サンプルは源から直接採取され得る。たとえば、血液または唾液が個体から直接採取され得る。代わりに、サンプルは源から直接採取されないことがある。次いで、1つまたは複数のプロセッサは、シーケンシングのためのサンプルを調製するようにシステムに指示する。この調製は、外来の物質を除去することおよび/または何らかの物質(たとえば、DNA)を隔離することを含み得る。生体サンプルは、特定のアッセイのための特徴を含むように調製され得る。たとえば、生体サンプルは、sequencing-by-synthesis(SBS)のために調製され得る。いくつかの実装形態では、調製することは、ゲノムのいくつかの領域の増幅を含み得る。たとえば、調製することは、STRおよび/またはSNRを含むことが知られている所定の遺伝子座を増幅することを含み得る。遺伝子座は、所定のプライマー配列を使用して増幅され得る。

10

【0278】

次に、1つまたは複数のプロセッサは、サンプルをシーケンシングするようにシステムに指示する。シーケンシングは、様々な既知のシーケンシングプロトコルを通じて実行され得る。特定の実装形態では、シーケンシングはSBSを含む。SBSでは、複数の蛍光ラベリングされたヌクレオチドが、光学基板の表面(たとえば、フローセルの中のチャンネルを少なくとも部分的に画定する表面)上に存在する増幅されたDNAの複数のクラスタ(場合によっては数百万個のクラスタ)をシーケンシングするために使用される。フローセルはシーケンシングのための核酸サンプルを含むことがあり、ここでフローセルは適切なフローセルホルダ内に配置される。

20

【0279】

核酸は、未知の標的配列に隣接する既知のプライマー配列を備えるように調製され得る。最初のSBSシーケンシングサイクルを開始するために、1つまたは複数の異なるようにラベリングされたヌクレオチド、およびDNAポリメラーゼなどが、流体サブシステムによってフローセルの中へと/フローセルを流れて流され得る。単一のタイプのヌクレオチドが一度に追加されるか、または、シーケンシング手順において使用されるヌクレオチドが反転可能な末端の性質を持つように特別に設計されるかのいずれかであってよく、これにより、シーケンシング反応の各サイクルが、いくつかのタイプのラベリングされたヌクレオチド(たとえば、A、C、T、G)の存在下で同時に発生することが可能になる。ヌクレオチドは、蛍光色素などの検出可能なラベル部分を含み得る。4つのヌクレオチドが一緒に混合される場合、ポリメラーゼは組み込むべき正しい塩基を選択することが可能であり、各配列は一塩基だけ延長される。組み込まれないヌクレオチドは、洗浄液をフローセルに流すことによって洗い落とされ得る。1つまたは複数のレーザーが、核酸を励起して蛍光を誘導し得る。核酸から放出される蛍光は組み込まれた塩基の蛍光色素に基づき、異なる蛍光色素は異なる波長の放出光を放出し得る。デブロッキング試薬が、延長され検出されたDNA鎖から反転可能な末端グループを除去するためにフローセルに追加され得る。次いでデブロッキング試薬が、洗浄液をフローセルに流すことによって洗い落とされ得る。そうすると、フローセルは、上に記載されたようなラベリングされたヌクレオチドの導入で開始するシーケンシングのさらなるサイクルの準備ができる。流体および検出の操作は、シーケンシングの実行を完了するために何回か繰り返され得る。例示的なシーケンシング方法は、たとえば、Bentley他、Nature 456:53-59(2008)、国際特許出願公開第WO 04/018497号、米国特許第7057026号、国際特許出願公開第WO 91/06678号、国際特許出願公開第WO 07/123744号、米国特許第7329492号、米国特許第7211414号、米国特許第7315019号、米国特許第7405281号、および米国特許出願公開第2008/0108082号において説明されており、これらの各々が参照によって本明細書において引用される。

30

40

【0280】

50

いくつかの実装形態では、核酸は表面に付着され、シーケンシングの前または間に増幅され得る。たとえば、増幅は、表面上に核酸クラスタを形成するためにブリッジ増幅を使用して行われ得る。有用なブリッジ増幅方法は、たとえば、米国特許第5641658号、米国特許出願公開第2002/0055100号、米国特許第7115400号、米国特許出願公開第2004/0096853号、米国特許出願公開第2004/0002090号、米国特許出願公開第2007/0128624号、および米国特許出願公開第2008/0009420号において説明されており、これらの各々の全体が参照によって本明細書において引用される。表面上で核酸を増幅するための別の有用な方法は、たとえば、Lizardi他、Nat. Genet. 19:225-232(1998)、および米国特許出願公開第2007/0099208A1号において説明されるようなローリングサークル増幅(RCA)であり、これらの各々が参照によって本明細書において引用される。

10

**【0281】**

1つの例示的なSBSプロトコルは、たとえば、国際特許出願公開第WO 04/018497号、米国特許出願公開第2007/0166705A1号、および米国特許第7057026号において説明されるような、除去可能な3'ブロックを有する修正されたヌクレオチドを利用し、これらの各々が参照によって本明細書において引用される。たとえば、SBS試薬の反復されるサイクルが、たとえばブリッジ増幅プロトコルの結果として、標的核酸が付着されたフローセルに導入され得る。核酸クラスタは、直線化溶液を使用して単鎖の形態へと変換され得る。直線化溶液は、たとえば、各クラスタの1本の鎖を開裂することが可能な制限エンドヌクレアーゼを含み得る。とりわけ化学開裂(たとえば、過ヨード酸塩を用いたジオール結合の開裂)、熱またはアルカリへの曝露によるエンドヌクレアーゼ(たとえば、米国マサチューセツ州イプスウィッチのNEBにより供給されるような「USER」、部品番号M5505S)を用いた開裂による無塩基サイトの開裂、そうされなければデオキシリボヌクレオチドからなる増幅産物へと組み込まれるリボヌクレオチドの開裂、光化学開裂またはペプチドリナーの開裂を含む、開裂の他の方法が、制限酵素またはニックング酵素に対する代替として使用され得る。直線化操作の後で、シーケンシングプライマーは、シーケンシングされるべき標的核酸へのシーケンシングプライマーのハイブリダイゼーションのための条件下で、フローセルに導入され得る。

20

**【0282】**

次いで、フローセルが、単一のヌクレオチドの追加によって各標的核酸にハイブリダイゼーションされるプライマーを延長するための条件下で、除去可能な3'ブロックおよび蛍光ラベルを伴う修正されたヌクレオチドを有するSBS延長試薬と接触させられ得る。単一のヌクレオチドだけが各プライマーに追加され、それは、修正されたヌクレオチドが、シーケンシングされているテンプレートの領域と相補的な成長中のポリヌクレオチド鎖へと組み込まれると、さらなる配列延長を指示するために利用可能な自由な3'-OH基がないので、ポリメラーゼがさらなるヌクレオチドを追加できないからである。SBS延長試薬は、除去され、放射線による励起のもとでサンプルを保護する構成要素を含む走査試薬により置き換えられ得る。走査試薬の例示的な構成要素は、米国特許出願公開第2008/0280773A1号および米国特許出願第13/018255号において説明され、これらの各々が参照によって本明細書に引用される。次いで、延長された核酸が、走査試薬の存在下で蛍光により検出され得る。蛍光が検出されると、3'ブロックが、使用されるブロッキンググループに適切なデブロック試薬を使用して除去され得る。それぞれのブロッキンググループに対して有用な例示的なデブロック試薬は、国際特許出願公開第WO004018497号、米国特許出願公開第2007/0166705A1号、および米国特許第7057026号において説明されており、これらの各々が参照によって本明細書において引用される。デブロック試薬は、3'OH基を有する延長されたプライマーにハイブリダイゼーションされる標的核酸を残して洗浄されてよく、このプライマーはこれで、さらなるヌクレオチドの追加が可能になる。したがって、延長試薬、走査試薬、およびデブロック試薬を追加するサイクルは、操作のうちの1つまたは複数の間の任意選択の洗浄とともに、所望の配列が得られるまで繰り返され得る。上記のサイクルは、修正されたヌクレオチドの各々に異なるラベルが付けられているとき、特定の塩基に対応することが知られている、サイクルごとに単一の延長試薬導入操作を使用して行わ

30

40

50

れ得る。異なるラベルが、各組み込み操作の間に追加されるヌクレオチドの区別を容易にする。代わりに、各サイクルは、延長試薬導入の別個の操作と、それに続く走査試薬導入と検出の別

個の操作とを含むことがあり、この場合、ヌクレオチドのうち2つ以上が同じラベルを有することが可能であり、それらを導入の既知の順序に基づいて区別することができる。

#### 【0283】

シーケンシング操作は特定のSBSプロトコルに関して上で論じられたが、シーケンシングのための他のプロトコルおよび様々な他の分子分析法のいずれもが、必要に応じて行われ得ることが理解されるであろう。

#### 【0284】

次いで、システムの1つまたは複数のプロセッサは、後続の分析のためのシーケンシングデータを受け取る。シーケンシングデータは、.BAMファイルなどの様々な方式でフォーマットされ得る。シーケンシングデータは、たとえばいくつかのサンプルリードを含み得る。シーケンシングデータは、ヌクレオチドの対応するサンプル配列を有する複数のサンプルリードを含み得る。1つだけのサンプルリードが論じられるが、シーケンシングデータは、たとえば、数百個、数千個、数十万個、または数百万個のサンプルリードを含み得ることを理解されたい。異なるサンプルリードは異なる数のヌクレオチドを有し得る。たとえば、サンプルリードは、10個のヌクレオチドから約500個以上のヌクレオチドにまでわたり得る。サンプルリードは源のゲノム全体にわたり得る。一例として、サンプルリードは、疑わしいSTRまたは疑わしいSNPを有する遺伝子座などの、所定の遺伝子座の方を向いている。

#### 【0285】

各サンプルリードは、サンプル配列、サンプルフラグメント、または標的配列と呼ばれ得る、ヌクレオチドの配列を含み得る。サンプル配列は、たとえば、プライマー配列、ランキング配列、および標的配列を含み得る。サンプル配列内のヌクレオチドの数は、30個、40個、50個、60個、70個、80個、90個、100個以上を含み得る。いくつかの実装形態では、サンプルリード(またはサンプル配列)のうち1つまたは複数は、少なくとも150個のヌクレオチド、200個のヌクレオチド、300個のヌクレオチド、400個のヌクレオチド、500個のヌクレオチド、またはそれより多くを含む。いくつかの実装形態では、サンプルリードは、1000個を超えるヌクレオチド、2000個を超えるヌクレオチド、またはそれより多くを含み得る。サンプルリード(またはサンプル配列)は、一端または両端にプライマー配列を含み得る。

#### 【0286】

次に、1つまたは複数のプロセッサは、シーケンシングデータを分析して、潜在的なバリエーションコールおよびサンプルバリエーションコールのサンプルバリエーション頻度を取得する。この操作は、バリエーションコールアプリケーションまたはバリエーションコーラとも呼ばれ得る。したがって、バリエーションコーラはバリエーションを特定または検出し、バリエーション分類器は検出されたバリエーションを体細胞性または生殖細胞性であるものとして分類する。代替的なバリエーションコーラが本明細書の実装形態に従って利用されることがあり、ここで、異なるバリエーションコーラは、実行されているシーケンシング操作のタイプ、対象のサンプルの特徴などに基づき使用され得る。バリエーションコールアプリケーションの1つの非限定的な例は、<https://github.com/Illumina/Pisces>においてホストされ、論説Dunn, TamsenおよびBerry, GwennおよびEmig-Agius, DorotheaおよびJiang, YuおよびIyer, AnitaおよびUdar, NitinおよびStromberg, Michael、(2017)、Pisces: An Accurate and Versatile Single Sample Somatic and Germline Variant Caller、595-595、10.1145/3107411.3108203において説明される、Illumina Inc.(カリフォルニア州サンディエゴ)によるPisces(商標)アプリケーションであり、上記の論説の完全な主題の全体が、参照によって本明細書において引用される。

#### 【0287】

そのようなバリエーションコールアプリケーションは4つの順番に実行されるモジュールを

10

20

30

40

50



備え得る。

【0288】

(1)Pisces Read Stithcer:BAMの中の対になっているリード(同じ分子のリード1およびリード2)をコンセンサスリードへとステッチングすることによってノイズを減らす。出力はステッチングされたBAMである。

【0289】

(2)Pisces Variant Caller:小さいSNV、挿入および欠失をコールする。Piscesは、リード境界、基本的なフィルタリングアルゴリズム、および単純なボワソンのバリエーション信頼性スコアリングアルゴリズムによって分解される、合祖バリエーションへのバリエーション折り畳み(variant-collapsing)アルゴリズムを含む。出力はVCFである。

10

【0290】

(3)Pisces Variant Quality Recalibrator(VQR):バリエーションコールが熱損傷またはFFPE脱アミノ化と関連付けられるパターンに圧倒的に従う場合、VQRステップは疑わしいバリエーションコールのバリエーションQスコアを下げる。出力は調整されたVCFである。

【0291】

(4)Pisces Variant Phaser(Scylla):クローンの垂集団から少数のバリエーションを複雑なアレルへと組み立てるために、read-backed greedy clustering法を使用する。このことは、下流のツールによる機能的な結果のより正確な決定を可能にする。出力は調整されたVCFである。

20

【0292】

加えて、または代わりに、この操作は、<https://github.com/Illumina/strelka>においてホストされ、論説T Saunders, ChristopherおよびWong, WendyおよびSwamy, SajaniおよびBecq, JenniferおよびJ Murray, LisaおよびCheetham, Keira、(2012)、Strelka: Accurate somatic small-variant calling from sequenced tumor-normal sample pairs、Bioinformatics (Oxford, England)、28、1811-7、10.1093/bioinformatics/bts271において説明される、Illumina Inc.によるバリエーションコールアプリケーションStrelka(商標)アプリケーションを利用することがあり、上記の論説の完全な主題の全体が、参照によって本明細書において明確に引用される。さらに、加えて、または代わりに、この操作は、<https://github.com/Illumina/strelka>においてホストされ、論説Kim, S、Scheffler, K、Halpern, A.L.、Bekritsky, M.A、Noh, E、Kallberg, M、Chen, X、Beyter, D、Krusche, P、およびSaunders, C.T、(2017)、Strelka2: Fast and accurate variant calling for clinical sequencing applicationsにおいて説明される、Illumina Inc.によるバリエーションコールアプリケーションStrelka2(商標)アプリケーションを利用することがあり、上記の文書の完全な主題の全体が、参照によって本明細書において明確に引用される。その上、加えて、または代わりに、この操作は、<https://github.com/Illumina/Nirvana/wiki>においてホストされ、論説Stromberg, MichaelおよびRoy, RajatおよびLajugie, JulienおよびJiang, YuおよびLi, HaochenおよびMargulies, Elliott、(2017)、Nirvana: Clinical Grade Variant Annotator、596-596、10.1145/3107411.3108204において説明される、Illumina Inc.によるNirvana(商標)アプリケーションなどのバリエーションアノテーション/コールツールを利用することがあり、上記の文書の完全な主題の全体が、参照によって本明細書において明確に引用される。

30

40

【0293】

そのようなバリエーションアノテーション/コールツールは、Nirvanaにおいて開示されるものなどの異なるアルゴリズム技法を適用することができる。

【0294】

a. 区間アレイ(Interval Array)を用いてすべての重複する転写産物を特定する。機能的なアノテーションのために、バリエーションと重複するすべての転写産物を特定することができ、区間木を使用することができる。しかしながら、区間のセットは静的であり得るので、区間木を区間アレイへとさらに最適化することが可能であった。区間木は、 $O(\min(n, k \lg n))$ 時間においてすべての重複する転写産物を返し、ここでnは木の中の区間の数であ

50

り、 $k$ は重複する区間の数である。実際には、 $k$ は大半のバリエーションに対して $n$ と比較して本当に小さいので、区間木上での実効的なランタイムは $O(k \lg n)$ である。我々は、最初の重複する区間を見つけて、次いで残りの $(k-1)$ にわたって数え上げるだけでよいように、ソートされたアレイにすべての区間が記憶されるような区間アレイを作成することによって、 $O(\lg n+k)$ へと改善した。

#### 【0295】

b. CNV/SV (Yu): コピー数変異および構造変異のためのアノテーションが提供され得る。小さいバリエーションのアノテーションと同様に、SVと重複する転写産物、および以前に報告された構造変異も、オンラインデータベースにおいてアノテートされ得る。小さいバリエーションと異なり、すべての重複する転写産物がアノテートされる必要はなく、それは、あまりにも多くの転写産物が大きなSVと重複するからである。代わりに、部分的な重複遺伝子に属するすべての重複する転写産物がアノテートされ得る。具体的には、これらの転写産物に対して、影響を受けるイントロン、エクソン、および構造変異により引き起こされる結果が報告され得る。すべての重複する転写産物の出力を許可するための選択肢が可能であるが、遺伝子名、転写産物との正規の (canonical) 重複であるか部分的な重複であるかのフラグなどの、これらの転写産物の基本情報が報告され得る。各SV/CNVに対して、これらのバリエーションが研究されているかどうか、および異なる集団におけるそれらの頻度を知ることに関心事である。したがって、1000 genomes、DGV、およびClinGenなどの外部データベースにおいて重複するSVを報告した。どのSVが重複しているかを決定するための恣意的なカットオフを使用するのを避けるために、代わりに、すべての重複する転写産物が使用されてよく、相互の重複率、すなわち、重複する長さをこれらの2つのSVの長さのうちの短い方で割ったものが計算されてよい。

10

20

#### 【0296】

c. 補足アノテーションを報告する。補足アノテーションには、小さいバリエーションと構造バリエーション (SV) という2つのタイプがある。SVは、区間としてモデル化されてよく、重複するSVを特定するために上で論じられた区間アレイを使用することができる。小さいバリエーションは、点としてモデル化され、場所および (任意選択で) アレルによって照合される。したがって、それらは二分探索様のアルゴリズムを使用して検索される。補足アノテーションデータベースは非常に大きいことがあるので、補足アノテーションが存在するファイル位置に染色体場所をマッピングするために、はるかにより小さなインデックスが作成される。インデックスは、場所を使用して二分探索され得るオブジェクト (染色体場所とファイル位置からなる) のソートされたアレイである。インデックスサイズを小さく保つために、(ある最大のカウンタまでの) 複数の場所が、第1の場所の値と後続の場所に対する差分のみを記憶する1つのオブジェクトへと圧縮される。二分探索を使用するので、ランタイムは $O(\lg n)$ であり、 $n$ はデータベースの中の項目の数である。

30

#### 【0297】

d. VEP キャッシュファイル

#### 【0298】

e. 転写産物データベース: 転写産物キャッシュ (キャッシュ) および補足データベース (SADB) ファイルは、転写産物および補足アノテーションなどのデータオブジェクトのシリアル化されたダンプである。キャッシュのためのデータソースとして、Ensembl VEP キャッシュを使用する。キャッシュを作成するために、すべての転写産物が区間アレイに挿入され、アレイの最終状態がキャッシュファイルに記憶される。したがって、アノテーションの間に、事前に計算された区間アレイをロードしてそれについて探索を実行するだけでよい。キャッシュはメモリへとロードされて探索は非常に高速 (上で説明された) であるため、重複する転写産物を見つけることは Nirvana においては非常に高速である (総ランタイムの1%未満であると鑑定されている?)。

40

#### 【0299】

f. 補足データベース: SADB のデータソースは補足材料のもとで列挙される。小さいバリエーションに対する SADB は、データベースの中の各オブジェクト (参照名および場所によって

50

特定される)がすべての関連する補足アノテーションを保持するように、すべてのデータソースのk-wayマージによって生み出される。データソースファイルを解析する間に遭遇する問題は、Nirvanaのホームページにおいて詳細に文書化されている。メモリ使用量を制限するために、SAインデックスのみがメモリにロードされる。このインデックスは、補足アノテーションのためのファイル位置の高速なルックアップを可能にする。しかしながら、データはディスクからフェッチされなければならないので、補足アノテーションを追加することは、Nirvanaの最大のボトルネックであると特定されている(総ランタイムの約30%であると鑑定されている)。

#### 【0300】

g. 結果および配列オントロジー(Consequence and Sequence Ontology):Nirvanaの機能的アノテーション(提供されるとき)は、Sequence Ontology(SO)(<http://www.sequenceontology.org/>)ガイドラインに従う。時として、現在のSOにおける問題を特定して、アノテーションの状態を改善するためにSOチームと協力する機会があった。

#### 【0301】

そのようなバリエーションツールは、前処理を含み得る。たとえば、Nirvanaは、ExAC、EVS、1000 Genomes project、dbSNP、ClinVar、Cosmic、DGV、およびClinGenのような、外部データソースからの多数のアノテーションを含んでいた。これらのデータベースを完全に利用するには、それらからの情報をサニタイジングしなければならない。我々は、様々なデータソースからの存在する様々な矛盾に対処するための異なる戦略を実施した。たとえば、同じ場所および代替のアレルに対する複数のdbSNPエントリがある場合、すべてのidをカンマで分けられたidのリストへと加える。同じアレルに対する異なるCAF値を伴う複数のエントリがある場合、第1のCAF値を使用する。矛盾するExACエントリとEVSエントリに対して、サンプルカウントの数を考慮し、よりサンプルカウントの高いエントリが使用される。1000 Genome Projectsでは、矛盾するアレルのアレル頻度を除去した。別の問題は不正確な情報である。我々は主に、1000 Genome Projectsからアレル頻度情報を抽出したが、GRCh38について、情報フィールドにおいて報告されているアレル頻度が利用可能ではない遺伝子型を伴うサンプルを除外していないことに気付き、これは、すべてのサンプルに対して利用可能ではないバリエーションに対する頻度の低下につながる。我々のアノテーションの正確さを保証するために、個体レベルの遺伝子型のすべてを使用して真のアレル頻度を計算する。知られているように、同じバリエーションは異なるアラインメントに基づく異なる表現を有し得る。すでに特定されたバリエーションについての情報を正確に報告できることを確実にするために、異なるリソースからのバリエーションを前処理してそれらを一貫した表現にしなければならない。すべての外部データソースに対して、基準アレルと代替アレルの両方における重複したヌクレオチドを除去するためにアレルをトリミングした。ClinVarについては、xmlファイルを直接解析し、すべてのバリエーションに対する5'アラインメントを実行した。この手法はしばしばvcfファイルにおいて使用される。異なるデータベースは同じ情報のセットを含み得る。不必要な重複を避けるために、一部の重複した情報を除去した。たとえば、1000 genomesにおけるこれらのバリエーションをより詳細な情報とともにすでに報告したので、1000 genome projectsをデータソースとして有するDGVの中のバリエーションを除去した。

#### 【0302】

少なくともいくつかの実装形態によれば、バリエーションコールアプリケーションは、低頻度バリエーション、生殖細胞系コーリングなどに対するコールを提供する。非限定的な例として、バリエーションコールアプリケーションは、腫瘍のみのサンプルおよび/または腫瘍-正常ペアサンプルに対して実行され得る。バリエーションコールアプリケーションは、一塩基変異(SNV)、多塩基変異(MNV)、インデルなどを探索し得る。バリエーションコールアプリケーションは、バリエーションを特定しながら、シーケンシングまたはサンプル調製エラーによる不一致をフィルタリングする。各バリエーションに対して、バリエーションコールは、基準配列、バリエーションの場所、および可能性のあるバリエーション配列(たとえば、AからCへのSNV、またはAGからAへの欠失)を特定する。バリエーションコールアプリケーションは、サンプル配列(また

10

20

30

40

50

はサンプルフラグメント)、基準配列/フラグメント、およびバリエーションコールを、バリエーションが存在することを示すものとして特定する。バリエーションコールアプリケーションは、生フラグメントを特定し、生フラグメントの指定、潜在的なバリエーションコールを検証する生フラグメントの数のカウント、支持バリエーションが発生した生フラグメント内での場所、および他の関連する情報を特定し得る。生フラグメントの非限定的な例には、duplex stitchedフラグメント、simplex stitchedフラグメント、duplex un-stitchedフラグメント、およびsimplex un-stitchedフラグメントがある。

#### 【0303】

バリエーションコールアプリケーションは、.VCFファイルまたは.GVCFファイルなどの様々なフォーマットでコールを出力し得る。単なる例として、バリエーションコールアプリケーションは、MiSeqReporterパイプラインに含まれ得る(たとえば、MiSeq(登録商標)シーケンサ装置で実装されるとき)。任意選択で、このアプリケーションは様々なワークフローを用いて実装され得る。分析は、所望の情報を得るために指定された方式でサンプルリードを分析する、単一のプロトコルまたはプロトコルの組合せを含み得る。

10

#### 【0304】

次いで、1つまたは複数のプロセッサは、潜在的なバリエーションコールに関連して妥当性確認操作を実行する。妥当性確認操作は、以下で説明されるように、品質スコア、および/または階層化された検定のヒエラルキーに基づき得る。妥当性確認操作が、潜在的なバリエーションコールを確認または実証するとき、妥当性確認操作は(バリエーションコールアプリケーションからの)バリエーションコール情報をサンプル報告生成器に渡す。代わりに、妥当性確認操作が潜在的なバリエーションコールを無効とするとき、または失格と判定するとき、妥当性確認操作は対応する指示(たとえば、否定的インジケータ、コールなしインジケータ、無効コールインジケータ)をサンプル報告生成器に渡す。妥当性確認操作はまた、バリエーションコールが正しいことまたは無効なコール指定が正しいことの信頼性の程度に関する信頼性スコアを渡し得る。

20

#### 【0305】

次に、1つまたは複数のプロセッサがサンプル報告を生成して記憶する。サンプル報告は、たとえば、サンプルに関する複数の遺伝子座に関する情報を含み得る。たとえば、遺伝子座の所定のセットの各遺伝子座に対して、サンプル報告は、遺伝子型コールを提供すること、遺伝子型コールを行えないことを示すこと、遺伝子型コールの確実さについての信頼性スコアを提供すること、または、1つまたは複数の遺伝子座に関するアッセイについての潜在的な問題を示すことのうちの少なくとも1つを行い得る。サンプル報告はまた、サンプルを提供した個体の性別を示し、かつ/またはサンプルが複数の源を含むことを示し得る。本明細書では、「サンプル報告」は、ある遺伝子座もしくは遺伝子座の所定のセットのデジタルデータ(たとえば、データファイル)および/または遺伝子座もしくは遺伝子座のセットの印刷された報告を含み得る。したがって、生成することまたは提供することは、データファイルを作成することおよび/もしくはサンプル報告を印刷すること、またはサンプル報告を表示することを含み得る。

30

#### 【0306】

サンプル報告は、バリエーションコールが決定されたが妥当性確認されなかったことを示すことがある。バリエーションコールが無効であると決定されるとき、サンプル報告は、バリエーションコールの妥当性を確認できないという決定の基礎に関する追加の情報を示し得る。たとえば、報告の中の追加の情報は、生フラグメントの記述と、生フラグメントがバリエーションコールを支持または否定する程度(たとえば、カウント)とを含み得る。加えて、または代わりに、報告の中の追加の情報は、本明細書で説明される実装形態に従って得られる品質スコアを含み得る。

40

#### 【0307】

##### [バリエーションコールアプリケーション]

本明細書で開示される実装形態は、潜在的なバリエーションコールを特定するためにシーケンシングデータを分析することを含む。バリエーションコールは、以前に実行されたシーケン

50

シング操作について記憶されたデータに対して実行され得る。加えて、または代わりに、バリエーションコーリングは、シーケンシング操作が実行されている間にリアルタイムで実行され得る。サンプルリードの各々が、対応する遺伝子座を割り当てられる。サンプルリードは、サンプルリードのヌクレオチドの配列、または言い換えると、サンプルリード内のヌクレオチドの順序(たとえば、A、C、G、T)に基づいて、対応する遺伝子座に割り当てられ得る。この分析に基づいて、サンプルリードは、特定の遺伝子座の潜在的なバリエーション/アレルを含むものとして指定され得る。サンプルリードは、遺伝子座の潜在的なバリエーション/アレルを含むものとして指定された他のサンプルリードとともに収集(または集約または貯蔵)され得る。割当て操作はコーリング操作とも呼ばれることがあり、コーリング操作において、サンプルリードは特定の遺伝子座/座と関連付けられる可能性があるものとして特定される。サンプルリードは、サンプルリードを他のサンプルリードから区別するヌクレオチドの1つまたは複数の識別配列(たとえば、プライマー配列)を位置特定するために分析され得る。より具体的には、識別配列は、特定の遺伝子座と関連付けられるものとしてサンプルリードを他のサンプルリードから特定し得る。

10

20

30

40

50

#### 【0308】

割当て操作は、識別配列の一連の $n$ 個のヌクレオチドが選択配列のうちの1つまたは複数と実質的に一致するかどうかを決定するために、識別配列の一連の $n$ 個のヌクレオチドを分析することを含み得る。特定の実装形態では、割当て操作は、サンプル配列の最初の $n$ 個のヌクレオチドが選択配列のうちの1つまたは複数と実質的に一致するかどうかを決定するために、サンプル配列の最初の $n$ 個のヌクレオチドを分析することを含み得る。数 $n$ は様々な値を有することがあり、この値はプロトコルへとプログラムされることがあり、またはユーザにより入力されることがある。たとえば、数 $n$ は、データベース内の最短の選択配列のヌクレオチドの数として定義され得る。数 $n$ は所定の数であり得る。所定の数は、たとえば、10、11、12、13、14、15、16、17、18、19、20、21、22、23、24、25、26、27、28、29、または30個のヌクレオチドであり得る。しかしながら、他の実装形態では、より少数または多数のヌクレオチドが使用され得る。数 $n$ はまた、システムのユーザなどの個人によって選択されてもよい。数 $n$ は1つまたは複数の条件に基づき得る。たとえば、数 $n$ は、データベース内の最短のプライマー配列のヌクレオチドの数または指定された数のうちの小さい方の数として定義され得る。いくつかの実装形態では、15未満のヌクレオチドであるあらゆるプライマー配列が例外として指定され得るように、15などの $n$ の最小値が使用され得る。

#### 【0309】

いくつかの場合、識別配列の一連の $n$ 個のヌクレオチドは、選択配列のヌクレオチドと正確に一致しないことがある。それでも、識別配列は、選択配列とほぼ同一である場合、選択配列と実質的に一致し得る。たとえば、識別配列の一連の $n$ 個のヌクレオチド(たとえば、最初の $n$ 個のヌクレオチド)が、不一致が指定された数(たとえば、3)を超えずに、かつ/またはシフトが指定された数(たとえば、2)を超えずに、選択配列と一致する場合、サンプルリードが遺伝子座に対してコールされ得る。各不一致またはシフトが、サンプルリードとプライマー配列との間の差としてカウントされ得るように、規則が確立され得る。差の数が指定された数未満である場合、サンプルリードは、対応する遺伝子座(すなわち、対応する遺伝子座に割り当てられる)に対してコールされ得る。いくつかの実装形態では、サンプルリードの識別配列と遺伝子座に関連付けられる選択配列との間の差の数に基づく、マッチングスコアが決定され得る。マッチングスコアが指定されたマッチング閾値に合格する場合、選択配列に対応する遺伝子座は、サンプルリードに対する潜在的な座として指定され得る。いくつかの実装形態では、サンプルリードが遺伝子座に対してコールされるかどうかを決定するために、後続の分析が実行され得る。

#### 【0310】

サンプルリードがデータベースの中の選択配列のうちの1つと実質的に一致する(すなわち、厳密に一致する、または上で説明されたようにほぼ一致する)場合、サンプルリードは、選択配列と相関する遺伝子座に割り当てられ、または指定される。これは、座コーリ

ングまたは予備的座コーリングと呼ばれることがあり、選択配列に相関する遺伝子座に対してサンプルリードがコールされる。しかしながら、上で論じられたように、サンプルリードは2つ以上の遺伝子座に対してコールされ得る。そのような実装形態では、潜在的な遺伝子座のうち1つだけに対するサンプルリードをコールするために、または割り当てるために、さらなる分析が実行され得る。いくつかの実装形態では、基準配列のデータベースと比較されるサンプルリードは、ペアエンドシーケンシングからの最初のリードである。ペアエンドシーケンシングを実行するとき、サンプルリードに相関する第2のリード(生フラグメントを表す)が得られる。割り当ての後で、割り当てられたリードについて実行される後続の分析は、割り当てられたリードに対してコールされた遺伝子座のタイプに基づき得る。

10

**【0311】**

次に、潜在的なバリエーションコールを特定するために、サンプルリードが分析される。とりわけ、分析の結果は、潜在的なバリエーションコール、サンプルバリエーション頻度、基準配列、および対象のゲノム配列内でのバリエーションが発生した場所を特定する。たとえば、遺伝子座がSNPを含むことが知られている場合、遺伝子座に対してコールされた割り当てられたリードは、割り当てられたリードのSNPを特定するための分析を経ることがある。遺伝子座が多型の反復的なDNA要素を含むことが知られている場合、サンプルリード内の多型の反復的なDNA要素を特定するために、または特徴付けるために、割り当てられるリードが分析され得る。いくつかの実装形態では、割り当てられるリードがSTR座およびSNP座と実質的に一致する場合、警告またはフラグがサンプルリードに割り当てられ得る。サンプルリードは、STR座とSNP座の両方として指定され得る。この分析は、割り当てられたリードの配列および/または長さを決定するために、割り当てプロトコルに従って割り当てられたリードをアラインメントすることを含み得る。アラインメントプロトコルは、全体が参照によって本明細書において引用される、2013年3月15日に出願された国際特許出願第PCT/US2013/030867号(公開番号第WO 2014/142831号)において説明される方法を含み得る。

20

**【0312】**

次いで、1つまたは複数のプロセッサは、支持バリエーションが生フラグメント内の対応する場所に存在するかどうかを決定するために、生フラグメントを分析する。様々なタイプの生フラグメントが特定され得る。たとえば、バリエーションコーラは、元のバリエーションコールを妥当性確認するバリエーションを示すあるタイプの生フラグメントを特定し得る。たとえば、そのタイプの生フラグメントは、duplex stitchedフラグメント、simplex stitchedフラグメント、duplex un-stitchedフラグメント、またはsimplex un-stitchedフラグメントを表し得る。任意選択で、前述の例の代わりに、またはそれに加えて、他の生フラグメントが特定され得る。各タイプの生フラグメントを特定することに関連して、バリエーションコーラはまた、支持バリエーションが発生した生フラグメント内での場所、ならびに、支持バリエーションを示した生フラグメントの数のカウントを特定する。たとえば、バリエーションコーラは、生フラグメントの10個のリードが、特定の場所Xにおいて支持バリエーションを有するduplex stitchedフラグメントを表すことが特定されたことを示すものを、出力し得る。バリエーションコーラはまた、生フラグメントの5個のリードが、特定の場所Yにおいて支持バリエーションを有するsimplex un-stitchedフラグメントを表すことが特定されたことを示すものを、出力し得る。バリエーションコーラはまた、基準配列に対応し、したがって対象のゲノム配列における潜在的なバリエーションコールを妥当性確認する証拠を提供する支持バリエーションを含まなかった、生フラグメントの数を出力し得る。

30

40

**【0313】**

次に、支持バリエーション、ならびに支持バリエーションが発生した場所を含む、生フラグメントのカウントが維持される。加えて、または代わりに、(サンプルリードまたはサンプルフラグメントの中の潜在的なバリエーションコールの場所に対する相対的な)対象の場所において支持バリエーションを含まなかった生フラグメントのカウントが維持され得る。加えて、または代わりに、基準配列に対応し潜在的なバリエーションコールを確認または確認しない、生フラグメントのカウントが維持され得る。潜在的なバリエーションコールを支持する生フラ

50

グメントのカウントおよびタイプ、生フラグメントの中の支持バリエーションの場所、潜在的なバリエーションコールを支持しない生フラグメントのカウントなどを含む、決定された情報が、バリエーションコール妥当性確認アプリケーションに出力される。

#### 【0314】

潜在的なバリエーションコールが特定されるとき、プロセスは、潜在的なバリエーションコール、バリエーション配列、バリエーション場所、およびそれらと関連付けられる基準配列を示すものを出力する。エラーはコールプロセスに誤ったバリエーションを特定させ得るので、バリエーションコールは「潜在的な」バリエーションを表すように指定される。本明細書の実装形態によれば、誤ったバリエーションまたは偽陽性を減らして除去するために、潜在的なバリエーションコールが分析される。加えて、または代わりに、プロセスは、サンプルリードと関連付けられる1つまたは複数の生フラグメントを分析し、生フラグメントと関連付けられる対応するバリエーションコールを出力する。

10

#### 【0315】

##### [ 良性訓練セットの生成 ]

数百万個のヒトゲノムおよびエクソンがシーケンシングされているが、それらの臨床上の応用は、疾患を引き起こす変異を良性の遺伝的変異から区別することの難しさにより限られたままである。ここで我々は、他の霊長類の種における一般的なミスセンスバリエーションが、ヒトにおいて大部分が臨床的に良性であることを実証し、病原性の変異が除去のプロセスによって系統的に特定されることを可能にする。6種のヒト以外の霊長類の種の集団シーケンシングからの数十万個の一般的なバリエーションを使用して、88%の正確さで稀な疾患の患者における病原性の変異を特定し、ゲノムワイド有意性(genome-wide significance)で知的障害における14個の新たな遺伝子候補の発見を可能にする、深層ニューラルネットワークを訓練した。追加の霊長類の種からの一般的な変異の目録を作成することで、数百万個の有意性が不確かなバリエーションに対する解釈が改善し、ヒトゲノムシーケンシングの臨床上の利用がさらに進む。

20

#### 【0316】

診断シーケンシングの臨床上の使用可能性は、ヒトの集団における稀な遺伝子バリエーションを解釈しそれらの疾患リスクに対する影響を推測することが難しいことにより、限られている。臨床的に有意な遺伝子バリエーションは、それらの健康に対する有害な影響により、集団において極めて稀である傾向があり、大半については、ヒトの健康に対する影響が決定されていない。臨床的な有意性が不確かであるこれらのバリエーションが多数あること、およびそれらが稀であることは、個人化された医療および集団全体の健康スクリーニングに対するシーケンシングの採用に対する手強い障壁となっている。

30

#### 【0317】

大半の浸透性のメンデル性の疾患は集団において非常に有病率が低いので、集団における高頻度でのバリエーションの観察は、良性の結果を支持する強い証拠である。多様なヒトの集団にわたって一般的な変異を評価することは、良性のバリエーションの目録を作成するための有効な戦略であるが、現生人類における一般的な変異の総数は、祖先の多様性の大部分が失われた我々の種の最近の歴史におけるボトルネック事象により、限られている。現生人類の集団の研究は、過去15000~65000年以内の10000人未満の個人という有効個体数( $N_e$ )からの顕著な膨張を示しており、一般的な多型のプールが小さいことは、このサイズの集団における変異の容量が限られていることに由来する。基準ゲノムの中の7000万個の潜在的なタンパク質を変化させるミスセンス置換のうち、全体で0.1%を超える集団アレル頻度を持つものは、概ね1000個のうちの1個しか存在しない。

40

#### 【0318】

現生人類の集団以外では、チンパンジーが次に近い現存する種を構成し、99.4%のアミノ酸配列相同性を共有する。ヒトとチンパンジーにおけるタンパク質コーディング配列の近い相同性は、チンパンジーのタンパク質コーディングバリエーションに対して作用する純化選択が、同一状態であるヒトの変異の健康に対する結果もモデル化し得ることを示唆する。

50

## 【0319】

中立的な多型がヒトの祖先の系統(約 $4N_e$ 世代)において持続する平均時間は、種の分岐時間(約600万年前)の一部であるので、自然に発生するチンパンジーの変異は、平衡選択により維持されるハプロタイプの稀な事例を除き、偶然を除いて大部分が重複しない変異空間に及ぶ。同一状態である多型が2つの種において同様に健康に影響する場合、チンパンジーの集団における高いアレル頻度でのバリエーションの存在は、ヒトにおける良性の結果を示すはずであり、その良性の結果が純化選択によって確立されている既知のバリエーションの目録を拡大する。

## 【0320】

[ 結果-他の霊長類における一般的なバリエーションはヒトにおいて大部分が良性である ]

Exome Aggregation Consortium(ExAC)およびGenome Aggregation Database(gnomAD)において収集された123136人のヒトを含む、集約されたエクソンデータが最近利用可能になったことで、アレル頻度スペクトラムにわたるミスセンス変異と同義変異に対する自然選択の影響を測ることが可能になった。コホートにおいて1回しか観察されない稀なシングルトンバリエーションは、変異率に対するトリヌクレオチドコンテキストの影響を調整した後の、de novo変異によって予測される、予想される2.2/1のミスセンス/同義比とよく一致する(図49A、図51、ならびに図52A、図52B、図52C、および図52D)が、より高いアレル頻度では、観察されるミスセンスバリエーションの数は、自然選択による有害な変異の一掃により減少する。アレル頻度の増大に伴うミスセンス/同義比の段階的な低下は、集団頻度が<0.1%であるミスセンスバリエーションのかなりの部分が、健康な個人において観察されるにもかかわらず軽度に有害な結果を有することと一致する。これらの発見は、0.1%~約1%より高いアレル頻度を伴うバリエーションを、平衡選択および創始者効果により引き起こされるよく記録されている少数の例外を除いて、浸透性の遺伝性疾患に対しては良性である可能性が高いものとして除去するという、診療室において広く行われている経験的な実践を支持するものである。

## 【0321】

我々は、24体の親類ではない個体のコホートにおいて2回以上サンプリングされた、一般的なチンパンジーバリエーションを特定した。これらのバリエーションの99.8%が一般のチンパンジー集団において一般的である(アレル頻度(AF)>0.1%)ことが推定され、これは、これらのバリエーションが純化選択のふるいにすでにかかけられていることを示す。我々は、複数の配列アラインメントにおける1対1のマッピングを欠いているバリエーションとともに、延長された主要組織適合遺伝子複合体領域を平衡選択の既知の領域として除いて、対応する同一状態のヒトバリエーションに対するヒトアレル頻度スペクトラムを調査した(図49B)。一般的なチンパンジーバリエーションと同一状態であるヒトバリエーションについて、ミスセンス/同義比はヒトアレル頻度スペクトラムにわたって概ね一定であり(カイ二乗(2)検定により $P>0.5$ )、これは、ヒトの集団における、一般的なチンパンジーバリエーションに対する負の選択がないことと、2つの種におけるミスセンスバリエーションに対する選択係数が合致していることと一致する。一般的なチンパンジーバリエーションと同一状態であるヒトバリエーションにおいて観察される低いミスセンス/同義比は、チンパンジーのより大きな有効個体数( $N_e$ 約73000)と一致し、これは軽度に有害な変異のより効率的な除去を可能にする。

## 【0322】

対照的に、シングルトンチンパンジーバリエーション(コホートにおいて1回しかサンプリングされない)について、一般的なアレル頻度においてミスセンス/同義比の大幅な低下が観察され( $P<5.8 \times 10^{-6}$ 、図49C)、これは、シングルトンチンパンジーミスセンスバリエーションの24%が、0.1%より高いアレル頻度ではヒトの集団における純化選択によってフィルタリングされるであろうことを示している。この枯渇は、チンパンジーシングルトンバリエーションの大部分が、その健康に対する有害な影響によりいずれの種においても一般的なアレル頻度に達することが妨げられた、稀な有害変異であることを示している。我々は、シングルトンバリエーションの69%だけが、一般のチンパンジー集団において一般的である(AF>0.1%)と推定する。

10

20

30

40

50



## 【0323】

次に、6種のヒト以外の霊長類の種のうちの少なくとも1種において観察される変異と同一状態であるヒトバリエントを特定した。6種の各々における変異は、大型類人猿ゲノムプロジェクト(チンパンジー、ボノボ、ゴリラ、およびオランウータン)から確認され、または、霊長類ゲノムプロジェクト(アカゲザル、マーモセット)から一塩基多型データベース(dbSNP)に提出され、シーケンシングされた個体の数が限られていること、および各種について観察されたミスセンス/同義比(補足テーブル1)が低いことに基づいて、大部分が一般的なバリエントを表す。チンパンジーと同様に、6種のヒト以外の霊長類の種からのバリエントのミスセンス/同義比は、少ない割合(チンパンジーにおいて0.1%未満のアレル頻度、および他の種ではシーケンシングされた個体がより少ないことによってより低いアレル頻度のもとで、約16%)の稀なバリエントが含まれることにより予想される一般的なアレル頻度におけるミスセンス変異の軽度の枯渇(図49D、図53、図54、および図55、補足データファイル1)を除き、ヒトアレル頻度スペクトラムにわたって概ね等しいことを発見した。これらの結果は、同一状態のミスセンスバリエントに対する選択係数が、ヒトの祖先の系統から約3500万年前に分岐したと推定される少なくとも新世界ザルまでの霊長類の系統内で合致していることを示唆する。

10

## 【0324】

観察された霊長類のバリエントと同一状態であるヒトミスセンスバリエントは、ClinVarデータベースにおける良性の結果に対して強くエンリッチされることを発見した。有意性が不確かであるバリエントおよび矛盾するアノテーションを伴うバリエントを除外した後で、少なくとも1種のヒト以外の霊長類の種において存在するClinVarバリエントは、平均で90%の事例で良性または良性である可能性が高いものとしてアノテートされ、それと比較して、ClinVarミスセンスバリエント全般では35%である( $P < 10^{-40}$ 、図49E)。霊長類バリエントに対するClinVarアノテーションの病原性は、精選バイアスを減らすために1%より高いアレル頻度を伴うヒトバリエントを除く健康なヒトの同様のサイズのコホートをサンプリングすることから観察されるもの(約95%が良性または良性である可能性が高い結果、 $P = 0.07$ )より、わずかに高い。

20

## 【0325】

ヒトの遺伝学の分野は、ヒト変異の臨床上的影響を推論するためにモデル生物に長い間依存してきたが、大半の遺伝的に扱いやすい動物モデルまでの進化的距離が長いことで、モデル生物についての発見がヒトに対してどの程度一般化可能であるかについての懸念が生まれている。我々は、4種の追加の哺乳類の種(ネズミ、ブタ、ヤギ、ウシ)と2種のより遠縁の脊椎動物(ニワトリ、ゼブラフィッシュ)からの概ね一般的な変異を含めるように、霊長類の系統を超えて分析を拡張した。我々は、dbSNPにおいてゲノムワイドの変異の確認が十分にとれている種を選択し、これらが概ね一般的なバリエントであることを、ミスセンス/同義比が2.2/1よりはるかに小さいことに基づいて確認した。我々の霊長類の分析とは対照的に、より遠縁の種における変異と同一状態であるヒトミスセンス変異は、一般的なアレル頻度では顕著に枯渇しており(図50A)、この枯渇の程度はより長い進化的距離において増大する(図50Bおよび補足テーブル2および3)。

30

## 【0326】

ヒトにおいて有害であるが、より遠縁の種では高いアレル頻度で耐えているミスセンス変異は、同一状態のミスセンス変異に対する選択の係数が、ヒトとより遠縁の種との間でかなり離れていることを示す。それでも、より遠縁の哺乳類におけるミスセンスバリエントの存在は良性の結果の確率を高め、それは、一般的なアレル頻度において自然選択により枯渇するミスセンスバリエントの割合が、ヒトミスセンスバリエント全般について観察される約50%の枯渇率より低いからである(図49A)。これらの結果と一致して、ネズミ、ブタ、ヤギ、およびウシにおいて観察されているClinVarミスセンスバリエントは、良性の結果または良性である可能性が高い結果をアノテートされる確率が73%であり、それと比較して、霊長類の変異に対しては90%であり( $P < 2 \times 10^{-8}$ 、図50C)、ClinVarデータベース全体に対しては35%であることを発見した。

40

50

## 【0327】

家畜化によるアーティファクトではなく進化的距離が選択係数の相違の主な原因であることを確認するために、我々は、広範囲の進化的距離にわたって、種内多型の代わりに近縁の種のペア間での固定された置換を使用して、分析を繰り返した(図50D、補足テーブル4、および補足データファイル2)。我々は、種間の固定された置換と同一状態であるヒトミスセンスバリエーションの枯渇率が、進化的な枝長とともに増大し、家畜化を受けた種と比較して野生種に対する識別可能な差がないことを発見した。これは、同一状態の固定されたミスセンス置換の数が分岐した系統において偶然により予想されるものよりも低かったことを発見した、ハエおよび酵母菌における以前の成果と一致している。

## 【0328】

[バリエーションの病原性分類のための深層学習ネットワーク]

開示される技術は、バリエーションの病原性分類のための深層学習ネットワークを提供する。临床上の応用に対するバリエーション分類の重要性は、教師あり機械学習を問題の対処のために使用する多くの試みを引き起こしてきたが、これらの努力は、訓練のために確信をもってラベリングされた良性のバリエーションおよび病原性のバリエーションを含む適切なサイズの真実データセット(truth dataset)がないことにより、妨げられている。

## 【0329】

専門家により精選されたバリエーションの既存のデータベースはゲノム全体を代表しておらず、ClinVarデータベースの中のバリエーションの約50%がわずかに200個の遺伝子(ヒトのタンパク質コーディング遺伝子の約1%)に由来する。その上、系統的な研究により、多くの専門家のアノテーションには証拠が疑わしいものがあることが特定されており、単一の患者においてのみ観察され得る希なバリエーションを解釈することの難しさを示している。専門家の解釈はますます厳密になっているが、分類のガイドラインは、大部分が合意された習慣に沿って策定されており、既存の傾向を強めるリスクがある。人による解釈のバイアスを減らすために、最近に分類器は、一般的なヒト多型または固定されたヒト-チンパンジーの置換に対して訓練されているが、これらの分類器も、人により精選されたデータベース上で訓練された以前の分類器の予測スコアを入力として使用している。これらの様々な方法の性能の客観的なベンチマーキングは、独立でバイアスのない真実データセットがなければ、達成が難しい。

## 【0330】

6種のヒト以外の霊長類(チンパンジー、ボノボ、ゴリラ、オランウータン、アカゲザル、およびマモセット)からの変異は、一般的なヒト変異と重複しない300000個を超える固有のミスセンスバリエーションに寄与し、大部分が純化選択のふるいにかげられた良性の結果の一般的なバリエーションを表し、機械学習の手法に利用可能な訓練データセットを大きく拡大する。平均すると、各霊長類の種は、ClinVarデータベースの全体より多くのバリエーション(2017年11月現在で、有意性が不確かなバリエーションおよび矛盾するアノテーションを伴うバリエーションを除いた後で、約42000個のミスセンスバリエーション)に寄与する。加えて、この内容は人の解釈によるバイアスがない。

## 【0331】

一般的なヒトバリエーション(AF>0.1%)および霊長類の変異(補足テーブル5(図58))を備えるデータセットを使用して、我々は新しい深層残差ネットワークPrimateAIを訓練した。PrimateAIは、対象のバリエーションの側にあるアミノ酸配列および他の種におけるオソロガス配列アラインメントを入力として取り込む(図2および図3)。人により加工された特徴量を利用する既存の分類器と異なり、我々の深層学習ネットワークは、元の配列から直接特徴量を抽出するように学習する。タンパク質構造についての情報を組み込むために、二次構造および溶媒接触性を配列だけから予測するように別々のネットワークを訓練し、次いでこれらを完全なモデルにおけるサブネットワークとして含めた(図5および図6)。結晶化に成功している少数のヒトタンパク質を仮定すれば、元の配列から構造を推論することには、不完全なタンパク質構造および機能ドメインのアノテーションによるバイアスが避けられるという利点がある。ネットワークの全体の深さは、含まれるタンパク質構造とともに

10

20

30

40

50

、およそ400000個の訓練可能なパラメータを備える36層の畳み込みであった。

【0332】

良性のラベルを伴うバリエーションのみを使用して分類器を訓練するために、我々は、所与の変異が集団において一般的なバリエーションとして観察される可能性が高いかどうかということとして、予測問題を形作った。いくつかの要因が高いアレル頻度でのバリエーションの観察の確率に影響し、我々はそれらのバリエーションの有害性だけに興味がある。他の要因には、変異率、シーケンシングカバレッジなどの技術的なアーティファクト、および遺伝子変換などの中立的な遺伝的浮動に影響する要因がある。

【0333】

我々は、これらの区別できない要因の各々を考慮して、良性の訓練セットの中の各バリエーションをExACデータベースからの123136個のエクソンに存在しなかったミスセンス変異と照合し、良性のバリエーションと照合された対照群とを区別するように深層学習ネットワークを訓練した(図24)。ラベリングされていないバリエーションの数は、ラベリングされた良性の訓練データセットのサイズを大きく超えるので、良性の訓練データセットと照合されたラベリングされていないバリエーションの異なるセットを各々使用して、8個のネットワークを並列に訓練し、コンセンサス予測を得た。

10

【0334】

元のアミノ酸配列のみを入力として使用して、深層学習ネットワークは、てんかん、自閉症、および知的障害における主要な疾患遺伝子である電位依存性ナトリウムチャンネルSCN2Aについて示されるように(図20)、有用なタンパク質機能ドメインにおける残基に高い病原性スコアを正確に割り当てている。SCN2Aの構造は、各々が6つの膜貫通ヘリックス(S1~S6)を含む4つの相同なリピートを備える。膜の脱分極により、正に荷電したS4膜貫通ヘリックスが膜の細胞外の側に向かって動き、S4-S5リンカーを介してS5/S6孔形成ドメインを開口させる。てんかん性脳症の早期の兆候と臨床的に関連付けられる、S4、S4-S5リンカー、およびS5ドメインにおける変異は、遺伝子において最高の病原性スコアを有するものとしてネットワークにより予測され、健康な集団におけるバリエーションに対して枯渇している(補足テーブル6)。我々はまた、ネットワークが、ドメイン内の重要なアミノ酸の場所を認識し、転写因子のDNA接触残基および酵素の触媒残基などの、これらの場所における変異に最高の病原性スコアを割り当てていることも発見した(図25A、図25B、図25C、および図26)。

20

30

【0335】

深層学習ネットワークがタンパク質の構造および機能についての洞察を元の配列からどれだけ導いているかをより理解するために、ネットワークの最初の3層からの訓練可能なパラメータを視覚化した。これらの層内で、Granthamスコアなどのアミノ酸距離の既存の測定結果に近い、異なるアミノ酸の重みと重みの間の相関を、ネットワークが学習することが観察された(図27)。これらの初期の層の出力はより後の層の入力になり、深層学習ネットワークがデータの次第により高次の表現を構築することを可能にする。

【0336】

訓練を保留された10000個の一般的な霊長類バリエーションを使用して、既存の分類アルゴリズムを用いたネットワークの性能を比較した。すべての新たに生じたヒトミスセンスバリエーションの約50%が一般的なアレル頻度で純化選択によってフィルタリングされるので(図49A)、変異率およびシーケンシングカバレッジによって10000個の一般的な霊長類バリエーションと照合された10000個のランダムに選択されたバリエーションのセットに対する各分類器の50パーセントイルスコアを求め、その閾値における各分類器の正確さを評価した(図21D、図28A、および補足データファイル4)。10000個の保留された一般的な霊長類バリエーションに良性の結果を割り当てることについて、我々の深層学習ネットワーク(91%の正確さ)は他の分類器の性能(次に良いモデルで80%の正確さ)を上回った。

40

【0337】

ヒト変異データのみを用いて訓練されたネットワークの正確さ(図21D)と比較した場合、既存の方法を超える改善の概ね半分は、深層学習ネットワークを使用することに由来し

50

、半分は訓練データセットを霊長類の変異で補強することに由来する。ある臨床シナリオにおける有意性が不確かなバリエーションの分類を検定するために、神経発達障害の患者vs健康な対照群において発生するde novo変異を区別する、深層学習ネットワークの能力を評価した。有病率により、神経発達障害は稀な遺伝子疾患の最大のカテゴリのうちの1つを構成しており、最近のトリオシーケンシング研究は、de novoミスセンス変異およびタンパク質切断変異の中心的な役割を示唆している。

#### 【0338】

各々確信をもってコールされている、Deciphering Developmental Disorders(DDD)コホートからの4293人の影響を受けている個人におけるde novoミスセンスバリエーションvs Simon's Simplex Collection(SSC)コホートにおける2517人の影響を受けていない兄弟からのde novoミスセンスバリエーションを分類し、ウィルコクソンの順位和検定を用いて2つの分布の間での予測スコアの差を評価した(図21Eおよび図29Aおよび図29B)。深層学習ネットワークは、このタスクについて他の分類器を明確に上回った( $P < 10^{-28}$ 、図21Fおよび図28B)。その上、保留された霊長類バリエーションデータセットと、DDD症例群vs対照群データセットとに対する、様々な分類器の性能が相関付けられ(スピアマン  $\rho = 0.57$ 、 $P < 0.01$ )、全く異なる源および方法を使用しているにもかかわらず、病原性の評価について2つのデータセットの間で良好な一致を示した(図30A)。

#### 【0339】

次に、同じ遺伝子内で良性変異vs病原性変異を分類することについての、深層学習ネットワークの正確さを推定することを試みる。DDD集団の大部分が、影響を受けている第一度近親者のいない、影響を受けている子供のインデックスケースを備えると仮定すると、分類器がde novo優性遺伝モードを持つ遺伝子の病原性を過剰評価することによって正確さを釣り上げていないことを示すのが重要である。我々は、タンパク質切断変異( $P < 0.05$ )だけから計算された、DDD研究において疾患との関連について名目上有意であった605個の遺伝子に分析を制約した。これらの遺伝子内で、de novoミスセンス変異は、予想と比較して3/1エンリッチされており(図22A)、約67%が病原性であることを示している。

#### 【0340】

深層学習ネットワークは、遺伝子の同じセット内で病原性のde novoバリエーションと良性のde novoバリエーションを区別することが可能であり( $P < 10^{-15}$ 、図22B)、他の方法の性能を大きく上回った(図22Cおよび図28C)。0.803以上というパイナリカットオフでは(図22Dおよび図30B)、症例群におけるde novoミスセンス変異の65%が病原性であるものとして深層学習ネットワークにより分類され、それと比較して、対照群においてはde novoミスセンス変異の14%が病原性であり、これは88%という分類の正確さに対応する(図22Eおよび図30C)。神経発達障害における頻繁な不完全な浸透性および変化する表現性を考慮すると、この数字は、対照群において部分的に浸透している病原性バリエーションが含まれていることにより、我々の分類器の正確さをおそらく過小評価している。

#### 【0341】

##### [新しい遺伝子候補の発見]

病原性ミスセンス変異を階層化するために0.803以上の閾値を適用することは、DDD患者におけるde novoミスセンス変異のエンリッチメントを、1.5-foldからタンパク質切断変異(2.5-fold)に近い2.2-foldへと増大させ、一方で、予想を超えてエンリッチされるバリエーションの総数の3分の1未満を捨てる。このことは、統計能力をかなり高め、元のDDD研究ではゲノムワイド有意性閾値にこれまで達していなかった知的障害における14個の追加の遺伝子候補の発見を可能にしている(テーブル1)。

#### 【0342】

##### [専門家による精選との比較]

ClinVarデータベースからの最近の専門家により精選されたバリエーションに対する様々な分類器の性能を調査したが、ClinVarデータセットに対する分類器の性能は、保留された霊長類バリエーションデータセットとも、DDD症例群vs対照群データセットとも強く相関していなかったことを発見した(それぞれ $P = 0.12$ および $P = 0.34$ )(図31Aおよび図31B)。我々は、

既存の分類器には専門家の精選によるバイアスがあるという仮説を立てており、人の経験則は正しい方向にある傾向にあるものの最適ではないことがある。1つの例は、ClinVarにおける病原性バリエーションと良性バリエーションとの間のGranthamスコアの平均の差であり、これは、605個の疾患関連遺伝子内での、DDD症例群vs対照群におけるde novoバリエーションの差の2倍である(テーブル2)。それと比べて、専門家による精選は、タンパク質構造を、特に、他の分子と相互作用することが可能になり得る表面に曝露されている残基の重要性を、十分に活用していないように見える。我々は、ClinVar病原性変異とDDD de novo変異の両方が、予測される溶媒に曝露される残基と関連付けられるが、良性のClinVarバリエーションと病原性のClinVarバリエーションとの間の溶媒接触性の差はDDD症例群vs対照群について見られる差の半分にすぎないことを観察した。これらの発見は、Granthamスコアおよび保存率などの、専門家にとって解釈がより簡単な要因を優先する確認バイアスを示唆するものである。人により精選されたデータベース上で訓練された機械学習分類器は、これらの傾向を強化することが予想される。

10

20

30

40

50

#### 【0343】

我々の結果は、系統的な霊長類集団のシーケンシングが、臨床上のゲノム解釈を現在制約している、数百万個の有意性が不確かなヒトバリエーションを分類するための有効な戦略であることを示唆している。保留された一般的な霊長類バリエーションと臨床上のバリエーションの両方に対する我々の深層学習ネットワークの正確さは、ネットワークを訓練するために使用される良性バリエーションの数とともに高まる(図23A)。その上、6種のヒト以外の霊長類の種の各々からのバリエーションについての訓練は、ネットワークの性能の向上に独立に寄与し、一方で、より遠縁の哺乳類からのバリエーションについての訓練はネットワークの性能に負の影響を与える(図23Bおよび図23C)。これらの結果は、一般的な霊長類バリエーションが、浸透性のメンデル性疾患に関してヒトにおいて大部分が良性である一方で、より遠縁の種における変異については同じことが言えないという主張を、支持している。

#### 【0344】

本研究において調査されるヒト以外の霊長類ゲノムの数は、シーケンシングされてきたヒトゲノムおよびエクソンの数と比較して少ないが、これらの追加の霊長類は、一般的な良性変異についての不相応な量の情報に寄与していることに留意することが重要である。ExACを用いたシミュレーションは、一般的なヒトバリエーション(>0.1%アレル頻度)の発見はわずかに数百の個体の後ですぐに停滞するが(図56)、数百万人までのさらなる健康な集団のシーケンシングが主に追加の希なバリエーションに寄与することを示している。アレル頻度に基づいて大部分が臨床的に良性であることが知られている一般的なバリエーションと異なり、健康な集団における希なバリエーションは、劣性の遺伝的疾患または浸透性が不完全である優性の遺伝的疾患を引き起こし得る。各霊長類の種は一般的なバリエーションの異なるプールを持つので、各種の数十の個体をシーケンシングすることは、霊長類の系統における良性ミスセンス変異の目録を系統的に作るのに有効な戦略である。実際に、本研究で調査された6種のヒト以外の霊長類の種からの134の個体が、ExAC研究からの123136人のヒトの4倍近く多くの一般的なミスセンス変異に寄与している(補足テーブル5(図58))。数百の個体に関する霊長類集団のシーケンシング研究は、野生の保護区域および動物園に住む比較的少数の親類ではない個体でも現実的であり得るので、野生の集団に対する外乱が最小限になり、これはヒト以外の霊長類の保全および倫理的な取り扱いの観点から重要である。

#### 【0345】

現生人類の集団は、大半のヒト以外の霊長類の種よりはるかに遺伝的な多様性が低く、チンパンジー、ゴリラ、およびテナガザルと比べて、個体ごとの一塩基バリエーションの数が概ね半分であり、オランウータンと比べて個体ごとのバリエーションが3分の1である。ヒト以外の霊長類の種の大半の遺伝的多様性のレベルは知られていないが、多数の現存するヒト以外の霊長類の種により、潜在的な良性のヒトミスセンスの場所の大半が少なくとも1つの霊長類の種における一般的なバリエーションによってカバーされる可能性が高いと推定することが可能になり、病原性バリエーションが除去のプロセスによって系統的に特定されることが可能になる(図23D)。シーケンシングされるこれらの種のサブセットのみでも、訓練デ

ータサイズを大きくすることで、機械学習を用いたミスセンスの結果のより正確な予測が可能になる。最終的に、我々の発見はミスセンス変異に注目しているが、この戦略は、バリエーションが同一状態であるかどうかを明確に決定するための十分なアラインメントがヒトゲノムと霊長類ゲノムとの間にある保存された制御領域においては特に、非コーディング変異の結果を推論することにも適用可能であり得る。

#### 【0346】

504種の既知のヒト以外の霊長類の種のうち、約60%が密猟および大規模な生息地の喪失により絶滅に瀕している。これらの種における個体数の減少と起こり得る絶滅は、遺伝的多様性における代わりのない損失となり、これらの固有の代わりのない種と我々自身の両方に利益をもたらすであろう、緊急を要する世界的な保全の努力に対する動機となっている。

10

#### 【0347】

##### [ データ生成およびアラインメント ]

アプリケーションの中の座標は、複数の配列アラインメントを使用してhg19にマッピングされる他の種におけるバリエーションに対する座標を含む、ヒトゲノムbuild UCSC hg19/GRCh37を参照する。タンパク質コーディングDNA配列に対する正規の転写産物および99種の脊椎動物ゲノムの複数の配列アラインメントおよび枝長が、UCSCゲノムブラウザからダウンロードされた。

#### 【0348】

Exome Aggregation Consortium(ExAC)/Genome Aggregation Database(gnomAD exomes)v 2.0から、ヒトエクソン多型データを取得した。24体のチンパンジー、13体のボノボ、27体のゴリラ、および10体のオランウータンに対する全体のゲノムシーケンシングデータおよび遺伝子型を備える、大型類人猿ゲノムシーケンシングプロジェクトからの霊長類変異データを取得した。チンパンジーおよびボノボの別の研究からの35体のチンパンジーからの変異も含めたが、バリエーションコーリング方法の違いにより、集団分析からはこれらを除外し、深層学習モデルの訓練にのみそれらを使用した。加えて、アカゲザルの個体16体およびマーモセットの個体9体が、これらの種に対する元のゲノムプロジェクトにおける変異を評価するために使用されたが、個体レベルの情報が利用可能ではなかった。アカゲザル、マーモセット、ブタ、ウシ、ヤギ、ネズミ、ニワトリ、およびゼブラフィッシュについての変異データをdbSNPから取得した。dbSNPは追加のオランウータンバリエーションも含んでおり、我々はそれを深層学習モデルの訓練にのみ使用した。それは、個体の遺伝子型情報が集団分析に利用可能ではなかったからである。平衡選択による影響を避けるために、集団分析のための延長された主要組織適合遺伝子複合体領域(chr6:28,477,797-33,448,354)内からのバリエーションも除外した。

20

30

#### 【0349】

ヒトタンパク質コーディング領域へのオーソログな1対1のマッピングを確実にし、偽遺伝子へのマッピングを防ぐために、99種の脊椎動物の多種アラインメントを使用した。バリエーションが基準方向/代替方向のいずれかで発生した場合、バリエーションを同一状態であるものとして受け入れた。バリエーションがヒトと他の種の両方において同じ予測されるタンパク質コーディング結果を有することを確実にするために、ミスセンスバリエーションと同義バリエーションの両方に対して、コドンの中の他の2つのヌクレオチドが種間で同一であることを要求した。分析に含まれる各種からの多型は補足データファイル1において列挙され、詳細な尺度は補足テーブル1に示されている。

40

#### 【0350】

4つのアレル頻度カテゴリーの各々に対して(図49A)、96個の潜在的なトリヌクレオチドコンテキストの各々における同義バリエーションとミスセンスバリエーションの予想される数を推定するために、および、変異率を訂正するために(図51および補足テーブル7、8(図59))、イントロン領域における変異を使用した。我々はまた、同一状態のCpGジヌクレオチドおよび非CpGジヌクレオチドバリエーションを別々に分析し、ミスセンス/同義比が両方のクラスに対してアレル頻度スペクトラムにわたって平坦であったことを検証した。これは、CpGバ

50

リアントと非CpGバリエーションの両方に対して、それらの変異率の違いが大きいのにもかかわらず、我々の分析が適用できることを示している(図52A、図52B、図52C、および図52D)。

【0351】

[他の種における多型と同一状態であるヒトミスセンスバリエーションの枯渇率]

他の種に存在するバリエーションがヒトにおいて一般的なアレル頻度(>0.1%)で耐えられるかどうかを評価するために、他の種における変異と同一状態であったヒトバリエーションを特定した。バリエーションの各々に対して、それらをヒト集団におけるそれらのアレル頻度に基づいて、4つのカテゴリ(シングルトン、シングルトンより多い~0.01%、0.01%~0.1%、>0.1%)のうちの1つに割り当て、稀(<0.1%)なバリエーションと一般的な(>0.1%)なバリエーションとの間でのミスセンス/同義比(MSR)の低下を推定した。一般的なヒトアレル頻度(>0.1%)での同一状態のミスセンスバリエーションの枯渇率は、ヒトにおける一般的なアレル頻度で自然選択により除去されるのに十分に有害な他の種からのバリエーションの割合を示す。

10

【0352】

【数40】

$$\% \text{depletion} = \frac{\text{MSR}_{\text{rare}} - \text{MSR}_{\text{comm}}}{\text{MSR}_{\text{rare}}}$$

【0353】

ミスセンス/同義比と枯渇の割合は、種ごとに計算され、図50Bおよび補足テーブル2に示される。加えて、チンパンジーの一般的なバリエーション(図49B)、チンパンジーのシングルトンバリエーション(図49C)、および哺乳類バリエーション(図50A)について、稀なバリエーションと一般的なバリエーションとの間でのミスセンス/同義比の差が有意であったかどうかを検定するために、2x2の分割表上で相同性のカイ二乗検定( $\chi^2$ )を実行した。

20

【0354】

シーケンシングは大型類人猿ゲノムプロジェクトからの限られた数の個体についてのみ実行されたので、一般的なチンパンジー集団において稀(<0.1%)または一般的な(>0.1%)であったサンプリングされたバリエーションの割合を推定するために、ExACからのヒトアレル頻度スペクトラムを使用した。ExACアレル頻度に基づいて24人のヒトのコホートをサンプリングし、このコホートにおいて一度、または一度より多く観察されたミスセンスバリエーションを特定した。一度より多く観察されたバリエーションは一般の集団において一般的な(>0.1%)である確率が99.8%であったが、コホートにおいて一度しか観察されなかったバリエーションは一般の集団において一般的である確率が69%であった。より遠縁の哺乳類におけるミスセンスバリエーションに対する観察される枯渇が、よく保存されておりしたがってより正確にアラインメントされている遺伝子の区別できない影響によるものではなかったことを検証するために、ヒトと比較した11種の霊長類および50種の哺乳類の複数の配列アラインメントにおいて50%を超える平均ヌクレオチド相同性を持つ遺伝子のみを限定して、上記の分析を繰り返した(補足テーブル3参照)。

30

【0355】

これは、結果に大きな影響を与えることなく、分析からヒトタンパク質コーディング遺伝子の約7%を取り除いた。加えて、結果がバリエーションコーリングまたは家畜化によるアーティファクト(dbSNPから選択された種の大半が家畜化されているので)の問題により影響を受けなかったことを確実にするために、種間多型の代わりに、近縁の種のペアからの固定された置換を使用して分析を繰り返した(図50D、補足テーブル4、および補足データファイル2)。

40

【0356】

[ヒト、霊長類、哺乳類、および他の脊椎動物に対する多型データのClinVar分析]

他の種と同一状態であるバリエーションの臨床上的影響を調査するために、矛盾する病原性のアノテーションを持っていたバリエーションまたは有意性が不確かなバリエーションとしてのみ

50

ラベリングされたバリエーションを除いて、ClinVarデータベースをダウンロードした。補足テーブル9に示されるフィルタリングステップの後で、合計で、病原性カテゴリの中の24853個のミスセンスバリエーションおよび良性カテゴリの中の17775個のミスセンスバリエーションがある。

#### 【0357】

ヒト、ヒト以外の霊長類、哺乳類、および他の脊椎動物における変異と同一状態であった、病原性ClinVarバリエーションおよび良性ClinVarバリエーションの数をカウントした。ヒトについては、ExACアレル頻度からサンプリングされた30人のヒトのコホートをシミュレートした。各種に対する良性バリエーションと病原性バリエーションの数が補足テーブル10に示されている。

10

#### 【0358】

[モデル訓練のための良性バリエーションとラベリングされていないバリエーションの生成]

機械学習のために、ヒトおよびヒト以外の霊長類からの大部分が一般的である良性ミスセンスバリエーションの良性訓練データセットを構築した。このデータセットは、一般的なヒトバリエーション(>0.1%のアレル頻度、83546個のバリエーション)、ならびにチンパンジー、ボノボ、ゴリラ、およびオランウータン、アカゲザル、およびマーモセットからのバリエーション(301690個の固有の霊長類バリエーション)を備える。各源が寄与する良性訓練バリエーションの数が補足テーブル5に示されている。

#### 【0359】

トリヌクレオチドコンテキスト、シーケンシングカバレッジ、およびそれらの種とヒトとの間のアラインメント可能性を考慮するために、照合されたラベリングされた良性バリエーションのセットとバリエーションのラベリングされていないセットとを区別するように、深層学習ネットワークを訓練した。ラベリングされていない訓練データセットを得るために、正規のコーディング領域におけるすべての潜在的なミスセンスバリエーションで開始した。ExACからの123136個のエクソンにおいて観察されたバリエーションと、開始コドンまたは終止コドンにおけるバリエーションとを除外した。合計で、68,258,623個のラベリングされていないミスセンスバリエーションが生成された。これは、シーケンシングカバレッジが悪い領域、および霊長類バリエーションに対する照合されたラベリングされていないバリエーションを選択するときにヒトゲノムと霊長類ゲノムとの間で1対1のアラインメントがなかった領域について修正するために、フィルタリングされた。

20

30

#### 【0360】

ラベリングされた良性バリエーションの同じセットと、ラベリングされていないバリエーションの8個のランダムにサンプリングされたセットとを使用する、8個のモデルを訓練し、それらの予測の平均をとることによって、コンセンサス予測を得た。妥当性確認および検定のために、10000個の霊長類バリエーションの2つのランダムにサンプリングされたセットも除外し、それらについては訓練を保留した(補足データファイル3)。これらのセットの各々に対して、トリヌクレオチドコンテキストによって照合された10000個のラベリングされていないバリエーションをサンプリングし、これらを、異なる分類アルゴリズム間で比較するときに各分類器の閾値を正規化するために使用した(補足データファイル4)。他の実装形態では、2個~500個にわたる、より少数のまたは追加のモデルがアンサンブルにおいて使用され得る。

40

#### 【0361】

一方は一般的なヒトバリエーションのみを用いて訓練され、もう一方は一般的なヒトバリエーションと霊長類バリエーションの両方を含む完全な良性のラベリングされたデータセットを用いた訓練された、深層学習ネットワークの2つのバージョンの分類の正確さを評価した。

#### 【0362】

[深層学習ネットワークのアーキテクチャ]

各バリエーションに対して、病原性予測ネットワークは、対象のバリエーションを中心とする長さ51のアミノ酸配列と、二次構造および溶媒接触性ネットワーク(図2および図3)の出力とを、中心の場所において置換されるミスセンスバリエーションとともに入力として取り込む。

50



11種の霊長類のための1つの場所頻度行列と、霊長類を除く50種の哺乳類のための1つの場所頻度行列と、霊長類と哺乳類を除く38種の脊椎動物のための1つの場所頻度行列とを含む、3つの長さ51の場所頻度行列が、99種の脊椎動物の複数の配列アラインメントから生成される。

#### 【0363】

二次構造深層学習ネットワークは、各アミノ酸の場所における3状態の二次構造、すなわちヘリックス(H)、シート(B)、およびコイル(C)を予測する(補足テーブル11)。溶媒接触性ネットワークは、各アミノ酸の場所における3状態の溶媒接触性、すなわち、埋もれている(buried)(B)、中間(intermediate)(I)、および露出している(exposed)(E)を予測する(補足テーブル12)。両方のネットワークが、入力としてフランキングアミノ酸配列のみを取り込み、Protein DataBankにおける既知の冗長ではない結晶構造からのラベルを使用して訓練された(補足テーブル13)。事前訓練された3状態二次構造ネットワークおよび3状態溶媒接触性ネットワークへの入力のために、やはり長さが51であり深さが20である、すべての99種の脊椎動物に対する複数の配列アラインメントから生成された単一の長さ場所頻度行列を使用した。Protein DataBankからの既知の結晶構造についてネットワークを事前訓練した後で、二次構造および溶媒モデルに対する最終的な2つの層が除去され、ネットワークの出力は病原性モデルの入力に直接接続された。3状態二次構造予測モデルについて達成される最良の検定の正確さは79.86%であった(補足テーブル14)。結晶構造を有していた約4000個のヒトタンパク質に対するDSSP(Define Secondary Structure of Proteins)とアノテートされた構造ラベルを使用するとき、予測される構造ラベルのみを使用するときとニューラルネットワークの予測を比較すると、大きな差はなかった(補足テーブル15)。

10

20

#### 【0364】

病原性予測のための我々の深層学習ネットワーク(PrimateAI)と、二次構造および溶媒接触性を予測するための深層学習ネットワークの両方が、残基ブロックのアーキテクチャを採用した。PrimateAIの詳細なアーキテクチャは、(図3)および補足テーブル16(図4A、図4B、および図4C)において説明されている。二次構造および溶媒接触性を予測するためのネットワークの詳細なアーキテクチャは、図6および補足テーブル11(図7Aおよび図7B)および12(図8Aおよび図8B)において説明されている。

#### 【0365】

[10000個の霊長類パリアントの保留された検定セットに対する分類器性能のベンチマーキング]

深層学習ネットワーク、ならびに、データベースdbNSFPから予測スコアを取得した他の20個のこれまでに公開されている分類器のベンチマークをとるために、検定データセットにおいて10000個の保留された霊長類パリアントを使用した。10000個の保留された霊長類パリアント検定セットに対する分類器の各々の性能も図28Aにおいて与えられる。異なる分類器は大きく変動するスコア分布を有していたので、各分類器に対する50パーセントイル閾値を特定するために、トリヌクレオチドコンテクストによって検定セットと照合された10000個のランダムに選択されたラベリングされていないパリアントを使用した。方法間での公平な比較を確実にするために、その分類器に対して50パーセントイルの閾値で良

30

40

#### 【0366】

分類器の各々に対して、50パーセントイル閾値を使用して良性であるものとして予測される保留された霊長類検定パリアントの割合も、図28Aおよび補足テーブル17(図34)に示されている。PrimateAIの性能は、パリアントの場所におけるアラインメントされた種の数に関してロバストであり、哺乳類からの十分な保存情報が利用可能である限り全般的に良好な性能であることも示し、これは大半のタンパク質コーディング配列について当てはまる(図57)。

#### 【0367】

50

## [ DDD研究からのde novoバリエーションの分析 ]

DDD研究からの公開されているde novoバリエーションと、SSC自閉症研究における健康な兄弟の対照群からのde novoバリエーションとを取得した。DDD研究はde novoバリエーションの信頼性レベルを提供しており、我々は、バリエーションコーリングエラーによる潜在的な偽陽性として、閾値が0.1未満であるバリエーションをDDDデータセットから除外した。一実装形態では、全体で、DDDの影響を受けている個人から3512個のミスセンスde novoバリエーションと、健康な対照群からの1208個のミスセンスde novoバリエーションがあった。99種の脊椎動物の複数配列アラインメントのためにUCSCによって使用された正規の転写産物アノテーションは、DDDにより使用される転写産物アノテーションとわずかに異なり、ミスセンスバリエーションの総数の小さな違いをもたらしている。DDDの影響を受けている個人におけるde novoミスセンスバリエーションと、自閉症研究からの影響を受けていない兄弟の対照群におけるde novoミスセンスバリエーションとを、この分類方法が区別する能力について評価した。各分類器に対して、2つの分布に対する予測スコア間の差のウィルコクソンの順位和検定からのP値を報告した(補足テーブル17(図34))。

10

## 【 0 3 6 8 】

同じ疾患遺伝子内での良性変異と病原性変異を様々な分類器が区別する際の正確さを高めるために、DDDコホートにおけるde novoタンパク質切断変異についてエンリッチされた( $P < 0.05$ 、ポワソン正確検定)、605個の遺伝子のサブセットに対して分析を繰り返した(補足テーブル18)。これらの605個の遺伝子内で、予想を超えるde novoミスセンス変異の3/1エンリッチメントに基づき、DDDデータセットの中のde novoバリエーションの3分の2が病原性であり、3分の1が良性であったと推定した。最小限の不完全な浸透と、健康な対照群におけるde novoミスセンス変異が良性であったことを仮定した。各分類器に対して、同じ数の良性の予測または病原性の予測を生み出した閾値を、これらのデータセットにおいて観察される経験的な割合として特定し、この閾値を、症例群vs対照群におけるde novo変異を各分類器が区別する際の正確さを推定するためのバイナリカットオフとして使用した。受信者動作特性曲線を構築するために、de novo DDDバリエーションの病原性の分類を真陽性のコールとして扱い、健康な対照群における病原性としてのde novoバリエーションの分類を偽陽性のコールとして扱った。DDDデータセットは3分の1の良性のde novoバリエーションを含むので、理論的に完璧な分類器に対する曲線下面積(AUC)は1より小さい。したがって、良性バリエーションと病原性バリエーションを完璧に分離する分類器は、DDD患者におけるde novoバリエーションの67%を真陽性として、DDD患者におけるde novoバリエーションの33%を偽陽性として、対照群におけるde novoバリエーションの100%を真陰性として分類し、0.837という最大の可能なAUCを生む(図29Aおよび図29Bおよび補足テーブル19(図35))。

20

30

## 【 0 3 6 9 】

## [ 新しい遺伝子候補の発見 ]

観察されるde novo変異の数をヌル変異モデルのもとで予想される数と比較することによって、遺伝子におけるde novo変異のエンリッチメントを検定した。DDD研究において実行されるエンリッチメント分析を繰り返し、PrimateAIスコアが0.803を超えるde novoミスセンス変異のみをカウントするときに新たにゲノムワイド有意である遺伝子を報告した。0.803を超えるPrimateAI閾値を満たすミスセンスバリエーションの割合(ゲノム全体で概ねすべての潜在的なミスセンス変異の5分の1)によって、de novoの損害を与えるミスセンス変異に対するゲノムワイド期待値を調整した。DDD研究ごとに、各遺伝子は4つの検定を必要とし、1つはタンパク質切断エンリッチメントを検定し、1つはタンパク質を変化させるde novo変異のエンリッチメントを検定し、両方が、DDDコホートだけのために、および神経発達トリオシーケンシングコホートのより大きなメタ分析のために検定される。タンパク質を変化させるde novo変異のエンリッチメントは、コーディング配列内のミスセンスde novo変異のクラスタリングの検定と、Fisherの方法によって組み合わせられた(補足テーブル20、21)。各遺伝子に対するP値が4つの検定の最小値から取られ、ゲノムワイド有意性が $P < 6.757 \times 10^{-7}$ として決定された( $\alpha = 0.05$ 、4つの検定を用いた18500個の遺伝子)。

40

## 【 0 3 7 0 】

50

### [ ClinVar分類の正確さ ]

既存の分類器の大半は、ClinVar上で訓練される分類器からの予測スコアを使用するなどして、ClinVarコンテンツ上で直接または間接的にのいずれかで訓練されるので、2017年以降に追加されたClinVarバリエーションのみを使用するように、ClinVarデータセットの分析を限定した。最近のClinVarバリエーションと他のデータベースとの間にはかなりの重複があったので、ExACにおいて一般的なアレル頻度(>0.1%)で見つかるバリエーション、または、HGMD(Human Gene Mutation Database)、LOVD(Leiden Open Variation Database)、またはUniProt(Universal Protein Resource)に存在するバリエーションを除去するために、さらにフィルタリングを行った。有意性が不確かであるものとしてだけアノテーションされたバリエーションおよび矛盾するアノテーションを伴うバリエーションを取り除いた後で、良性のアノテーションを伴う177個のミスセンスバリエーションおよび病原性のアノテーションを伴う969個のミスセンスバリエーションが残った。これらのClinVarバリエーションを、深層学習ネットワークと他の分類方法の両方を使用してスコアリングした。各分類器に対して、同じ数の良性予測と病原性予測を生み出した閾値を、これらのデータベースにおいて観察される経験的な割合として特定し、この閾値を、各分類器の正確さを推定するためのバイナリカットオフとして使用した(図31Aおよび図31B)。

10

### 【 0 3 7 1 】

[ 訓練データサイズを大きくすることおよび訓練データの異なる源を使用することの影響 ]

深層学習ネットワークの性能に対する訓練データサイズの影響を評価するために、385236個の霊長類および一般的なヒトのバリエーションの良性とラベリングされた訓練セットから、バリエーションのサブセットをランダムにサンプリングし、背後の深層学習ネットワークアーキテクチャを同一に保った。各々の個別の霊長類の種からのバリエーションが分類の正確さに寄与する一方で、各々の個別の哺乳類の種からのバリエーションはより低い分類の正確さに寄与することを示すために、一実装形態に従って、83546個のヒトバリエーションと、各種に対するランダムに選択された一定の数のバリエーションとを備える訓練データセットを使用して、深層学習ネットワークを訓練し、背後のネットワークアーキテクチャを再び同じに保った。訓練セットに追加したバリエーションの一定の数(23380)は、ミスセンスバリエーションの数が最小である種、すなわちボノボにおいて利用可能なバリエーションの総数であった。各分類器に対する性能の中央値を得るために、訓練手順を5回繰り返した。

20

### 【 0 3 7 2 】

[ シーケンシングされる霊長類集団の数の増大に伴うすべての潜在的なヒトミスセンス変異の飽和 ]

ExACにおいて観察される一般的なヒトミスセンスバリエーション(>0.1%のアレル頻度)のトリヌクレオチドコンテキストに基づいてバリエーションをシミュレートすることによって、504種の現存する霊長類の種において存在する一般的なバリエーションによる、すべての約7000万個の潜在的なヒトミスセンス変異の予想される飽和を調査した。各霊長類の種に対して、ヒトにおいて観察される一般的なミスセンスバリエーションの数(アレル頻度が0.1%を超える約83500個のミスセンスバリエーション)の4倍をシミュレートした。それは、ヒトが、他の霊長類の種と比べて個体あたりのバリエーションの数が概ね半分であり、0.1%を超えるアレル頻度では、純化選択によって約50%のヒトミスセンスバリエーションが取り除かれているから

40

### 【 0 3 7 3 】

調査されるヒトのコホートのサイズの増大に伴って発見される一般的なヒトミスセンスバリエーション(>0.1%のアレル頻度)の割合をモデル化するために(図56)、ExACアレル頻度に従って遺伝子型をサンプリングし、これらのシミュレートされるコホートにおいて少なくとも1回観察された一般的なバリエーションの割合を報告した。

### 【 0 3 7 4 】

一実装形態では、PrimateAIスコアの現実的な応用のために、優性遺伝モードを伴う遺伝子においては、対照群と比較した症例群におけるde novoバリエーションのエンリッチメントに基づいて(図21D)、>0.8という閾値が病原性の可能性が高いという分類に対して好ま

50

しく、<0.6が良性である可能性が高いという分類に対するものであり、0.6~0.8が中間であり、劣性遺伝モードを伴う遺伝子においては、>0.7という閾値が病原性である可能性が高いという分類に対するものであり、<0.5が良性である可能性が高いという分類に対するものである。

#### 【0375】

図2は、本明細書で「PrimateAI」と呼ばれる、病原性予測のための深層残差ネットワークの例示的なアーキテクチャを示す。図2において、1Dは1次元畳み込み層を指す。予測される病原性は、0(良性)から1(病原性)までの目盛り上にある。ネットワークは、ヒトアミノ酸(AA)基準およびパリアントを中心とする代替配列(51個のAA)、99種の脊椎動物の種から計算された位置特定の重み行列(PWM)保存プロファイル、ならびに二次構造および溶媒接触性予測深層学習ネットワークの出力を入力として取り込み、この深層学習ネットワークは、3状態のタンパク質二次構造(ヘリックス-H、シート-B、およびコイル-C)と、3状態の溶媒接触性(埋もれている-B、中間-I、および露出している-E)とを予測する。

10

#### 【0376】

図3は、病原性分類のための深層学習ネットワークアーキテクチャであるPrimateAIの概略図を示す。モデルへの入力、基準配列と置換されるパリアントを伴う配列との両方に対するランキング配列の51個のアミノ酸(AA)と、霊長類、哺乳類、および脊椎動物のアラインメントからの3つの長さ51AAの位置特定の重み行列により表される保存率と、事前訓練された二次構造ネットワークおよび溶媒接触性ネットワークの出力(やはり長さは51A Aである)とを含む。

20

#### 【0377】

図4A、図4B、および図4Cは、病原性予測深層学習モデルPrimateAIの例示的なモデルアーキテクチャの詳細を示す、補足テーブル16である。形状はモデルの各層における出力テンソルの形状を指定し、活性化は層のニューロンに与えられる活性化である。モデルへの入力、パリアントの周りのランキングアミノ酸配列に対する位置特定の頻度行列(長さ51AA、深さ20)、ワンホット符号化された(one-hot encoded)ヒト基準配列および代替配列(長さ51AA、深さ20)、ならびに、二次構造および溶媒接触性モデルからの出力(長さ51A A、深さ40)である。

#### 【0378】

示される例は1D畳み込みを使用する。他の実装形態では、このモデルは、2D畳み込み、3D畳み込み、拡張畳み込みまたは膨張畳み込み、転置畳み込み(transposed convolution)、分離可能畳み込み(separable convolution)、および深さごとの(depthwise)分離可能畳み込みなどの、異なるタイプの畳み込みを使用することができる。一部の層は、シグモイドまたは双曲線正接などの飽和する非線形性と比較して確率的勾配降下の収束を大きく加速する、ReLU活性化関数も使用する。開示される技術によって使用され得る活性化関数の他の例には、parametric ReLU、leaky ReLU、および指数関数的線形ユニット(ELU)がある。

30

#### 【0379】

一部の層はバッチ正規化(IoffeおよびSzegedy、2015年)も使用する。バッチ正規化に関して、畳み込みニューラルネットワーク(CNN)における各層の分布は訓練の間に変化し、層によって変化する。このことは、最適化アルゴリズムの収束速度を低下させる。バッチ正規化はこの問題を克服するための技法である。バッチ正規化層の入力をxで表し、その出力をzを使用して表すと、バッチ正規化はxについての以下の変換を適用する。

40

#### 【0380】

#### 【数41】

$$z = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \gamma + \beta$$

50

## 【0381】

バッチ正規化は、 $\mu$ および  $\sigma$  を使用して入力 $x$ に対する平均-分散の正規化を適用し、 $\mu$ および  $\sigma$  を使用してそれを線形にスケールしてシフトする。正規化パラメータ  $\mu$  および  $\sigma$  は、指数移動平均と呼ばれる方法を使用して訓練セットにわたって現在の層に対して計算される。言い換えると、それらは訓練可能なパラメータではない。対照的に、 $\mu$  および  $\sigma$  は訓練可能なパラメータである。訓練の間に計算される  $\mu$  および  $\sigma$  の値は、推論の間にフォワードパスにおいて使用される。

## 【0382】

図5および図6は、タンパク質の二次構造および溶媒接触性を予測するために使用される深層学習ネットワークアーキテクチャを示す。モデルに対する入力は、RaptorXソフトウェア(Protein Data Bank配列について訓練するための)または99種の脊椎動物のアラインメント(ヒトタンパク質配列についての訓練および推論のための)によって生成される保存率を使用した、位置特定の重み行列である。長さが51AAである、第2の層から最後の層の出力は、病原性分類のための深層学習ネットワークに対する入力になる。

10

## 【0383】

図7Aおよび図7Bは、3状態二次構造予測深層学習(DL)モデルの例示的なモデルアーキテクチャの詳細を示す補足テーブル11である。形状はモデルの各層における出力テンソルの形状を指定し、活性化は層のニューロンに与えられる活性化である。モデルへの入力は、バリエーションの周りのフランキンガミノ酸配列に対する位置特定の頻度行列(長さ51AA、深さ20)であった。

20

## 【0384】

図8Aおよび図8Bは、3状態溶媒接触性予測深層学習モデルの例示的なモデルアーキテクチャの詳細を示す補足テーブル12である。形状はモデルの各層における出力テンソルの形状を指定し、活性化は層のニューロンに与えられる活性化である。モデルへの入力は、バリエーションの周りのフランキンガミノ酸配列に対する位置特定の頻度行列(長さ51AA、深さ20)であった。

## 【0385】

図20は、重要な機能ドメインに対してアノテートされた、SCN2A遺伝子における各アミノ酸の位置における予測される病原性スコアを示す。遺伝子に沿ってプロットされているのは、各アミノ酸の位置におけるミスセンス置換に対する平均PrimateAIスコアである。

30

## 【0386】

図21Dは、訓練を保留された10000個の一般的な霊長類バリエーションの検定セットに対する良性的結果を予測することにおける分類器の比較を示す。y軸は、各分類器の閾値を変異率について照合された10000個のランダムなバリエーションのセットについての50パーセントイルスコアへと正規化した後の、良性的であるものとして正しく分類された霊長類バリエーションの百分率を表す。

## 【0387】

図21Eは、Deciphering Developmental Disorders(DDD)患者において発生するde novoミスセンスバリエーションに対するPrimateAI予測スコアの分布を、影響を受けていない兄弟と比較して、対応するウィルコクソンの順位和のP値とともに示す。

40

## 【0388】

図21Fは、DDD症例群vs対照群におけるde novoミスセンスバリエーションを分離することにおける分類器の比較を示す。ウィルコクソンの順位和検定のP値が各分類器に対して示されている。

## 【0389】

図22A、図22B、図22C、図22D、および図22Eは、 $P < 0.05$ である605個のDDD遺伝子内の分類の正確さを示す。図22Aは、de novoタンパク質切断変異( $P < 0.05$ )に対して有意であった605個の関連する遺伝子内の、DDDコホートからの影響を受けている個人における予想を超えるde novoミスセンス変異のエンリッチメントを示す。図22Bは、605個の関連する遺伝子内での、DDD患者vs影響を受けていない兄弟において発生するde novoミスセンスバリア

50

ントに対するPrimateAI予測スコアの分布を、対応するウィルコクソンの順位和のP値とともに示す。

【0390】

図22Cは、605個の遺伝子内での、症例群vs対照群におけるde novoミスセンスバリエントを分離する際の様々な分類器の比較を示す。y軸は、各分類器に対するウィルコクソンの順位和検定のP値を示す。

【0391】

図22Dは、受信者動作特性曲線上で示される、様々な分類器の比較を、各分類器に対して示されるAUCとともに示す。

【0392】

図22Eは、各分類器に対する分類の正確さとAUCを示す。示される分類の正確さは、図22Aにおいて示されるエンリッチメントに基づいて予測されるのと同じ数の病原性バリエントと良性バリエントを分類器が予測するような閾値を使用した、真陽性と真陰性のエラー率の平均である。DDD de novoミスセンスバリエントの33%がバックグラウンドとなるという事実を考慮するために、完璧な分類器に対する最大の達成可能なAUCが点線で示される。

【0393】

図23A、図23B、図23C、および図23Dは、訓練のために使用されるデータの、分類の正確さに対する影響を示す。深層学習ネットワークは、完全なデータセット(385236個のバリエント)まで、増大する数の霊長類およびヒトの一般的なバリエントを用いて訓練される。図23Aにおいて、ネットワークの各々の分類性能のベンチマークが、10000個の保留された霊長類バリエントと、DDD症例群vs対照群におけるde novoバリエントとに対して正確さについてとられる。

【0394】

図23Bおよび図23Cは、一実装形態による、83546個の一般的なヒトバリエントと、単一の霊長類の種または哺乳類の種からの23380個のバリエントとを備えるデータセットを使用して訓練された、ネットワークの性能を示す。10000個の保留された霊長類バリエント(図23B)と、DDD症例群vs対照群におけるde novoミスセンスバリエント(図23C)とについてベンチマークがとられた、一般的な変異の異なる源を用いて訓練された各ネットワークに対して、結果が示されている。

【0395】

図23Dは、504種の現存する霊長類の種における同一状態の一般的なバリエント(>0.1%)による、すべての潜在的な良性のヒトミスセンスの位置の予想される飽和を示す。y軸は少なくとも1つの霊長類の種において観察されるヒトミスセンスバリエントの割合を示し、CpGミスセンスバリエントが緑で示され、すべてのミスセンスバリエントが青で示されている。各霊長類の種における一般的なバリエントをシミュレートするために、置き換えを伴うすべての潜在的な一塩基置換のセットからサンプリングし、ExACにおける一般的なヒトバリエント(>0.1%のアレル頻度)について観察されるトリヌクレオチドコンテキスト分布を照合した。

【0396】

図24は、一般的な霊長類バリエントの確認に対するシーケンシングカバレッジの影響を訂正することを示す。ヒト以外の霊長類の種において所与のバリエントを観察する確率は、ExAC/gnomADエクソンデータセットの中のその位置におけるシーケンシング深度と逆の相関がある。対照的に、より低いgnomADリードの深さは、その位置における一般的なヒトバリエント(>0.1%のアレル頻度)を観察する確率に影響を与えなかった。それは、シーケンシングされる多数のヒトエクソンが、一般的な変異の確認をほとんど保証されたものにするからである。ネットワークを訓練するための霊長類バリエントの各々に対する一致したバリエントを選ぶとき、あるバリエントを選ぶ確率は、変異率および遺伝子変換を考慮するためのトリヌクレオチドコンテキストに対する照合に加えて、シーケンシングの深さの影響について調整された。

10

20

30

40

50

## 【0397】

図25A、図25B、図25C、および図26は、開示されたニューラルネットワークによるタンパク質モチーフの認識を示す。図25A、図25B、および図25Cに関して、タンパク質ドメインのニューラルネットワークによる認識を例示するために、3つの異なるタンパク質ドメインの中の各アミノ酸位置におけるバリエーションに対する平均PrimateAIスコアを示す。図25Aにおいて、反復するGXXモチーフの中にグリシンを伴う、COL1A2というコラーゲン鎖が強調されている。コラーゲン遺伝子における臨床的に特定されている変異は、大部分がGXXリピートにおけるグリシンのミスセンス変異によるものであり、それは、これらがコラーゲンの正常な組み立てと干渉し、強いドミナントネガティブ効果を及ぼすからである。図25Bにおいて、IDSスルファターゼ酵素の活性サイトが強調されており、これは、ホルミルグリシンへと翻訳後修飾されるシステインを活性サイトに含む。図25Cにおいて、MYC転写因子のbHLHzipドメインが示されている。基本ドメインは、負に荷電した糖リン酸の骨格と相互作用する正に荷電したアルギニンおよびリジン残基(強調されている)を介して、DNAに接触する。ロイシンジッパードメインは、二量体化のために決定的に重要である、7個のアミノ酸だけ離隔されたロイシン残基(強調されている)を備える。

10

## 【0398】

図26は、バリエーションに対する予測される深層学習スコアへの、バリエーションの中および周りの各位置を摂動させることの影響を示す線のプロットを含む。バリエーションの周りの近くのアミノ酸(位置-25~+25)における入力を系統的に0に設定し、ニューラルネットワークによるバリエーションの予測される病原性の変化を測定した。プロットは、5000個のランダムに選択されたバリエーションに対する各々の近くのアミノ酸位置における摂動に対する、予測される病原性スコアの平均の変化を示す。

20

## 【0399】

図27は、重みの相関パターンがBLOSUM62スコア行列およびGranthamスコア行列を模倣することを示す。二次構造深層学習ネットワークの最初の3つの層からの重みの相関パターンは、BLOSUM62スコア行列およびGranthamスコア行列に類似するアミノ酸間の相関を示す。左のヒートマップは、ワンホット表現を使用して符号化されたアミノ酸間の二次構造深層学習ネットワークの2つの初期アップサンプリング層に続く、第1の畳み込み層からのパラメータ重みの相関を示す。中間のヒートマップは、アミノ酸のペア間のBLOSUM62スコアを示す。右のヒートマップはアミノ酸間のGrantham距離を示す。深層学習重みとBLOSUM62スコアとの間のPearson相関は $0.63(P=3.55 \times 10^{-9})$ である。深層学習重みとGranthamスコアとの間の相関は $-0.59(P=4.36 \times 10^{-8})$ である。BLOSUM62スコアとGranthamスコアとの間の相関は $-0.72(P=8.09 \times 10^{-13})$ である。

30

## 【0400】

図28A、図28B、および図28Cは、深層学習ネットワークPrimateAIと他の分類器の性能評価を示す。図28Aは、訓練を保留された10000個の霊長類バリエーションの検定セットに対する良性的結果を予測することにおける深層学習ネットワークPrimateAIの正確さと、SIFT、PolyPhen-2、CADD、REVEL、M-CAP、LRT、MutationTaster、MutationAssessor、FATHMM、PROVEAN、VEST3、MetaSVM、MetaLR、MutPred、DANN、FATHMM-MKL\_coding、Eigen、GenoCanyon、integrated\_fitCons、およびGERPを含む、他の分類器との比較とを示す。y軸は、変異率および遺伝子変換を考慮するためにトリヌクレオチドコンテキストについて霊長類バリエーションと照合された、10000個のランダムに選択されたバリエーションのセットを使用して、各分類器に対する閾値を50パーセントイルスコアへと正規化したことに基づいて、良性的であるものとして分類された霊長類バリエーションの百分率を表す。

40

## 【0401】

図28Bは、上で列挙された20個の既存の方法とともに、DDD症例群vs対照群におけるde novoミスセンスバリエーションを分離する際のPrimateAIネットワークの性能の比較を示す。y軸は、各分類器に対するウィルコクソンの順位和検定のP値を示す。

## 【0402】

図28Cは、上で列挙された20個の方法とともに、605個の疾患関連遺伝子内での、DDD症

50

例群vs影響を受けていない対照群におけるde novoミスセンスバリエントを分離する際のPrimateAIネットワークの性能の比較を示す。y軸は、各分類器に対するウィルコクソンの順位和検定のP値を示す。

【0403】

図29Aおよび図29Bは、4つの分類器の予測スコアの分布を示す。DDD症例群vs影響を受けていない対照群において発生するde novoミスセンスバリエントに対する、SIFT、PolyPhen-2、CADD、およびREVELを含む4つの分類器の予測スコアのヒストグラムを、対応するウィルコクソンの順位和のP値とともに示す。

【0404】

図30A、図30B、および図30Cは、605個の疾患関連遺伝子における病原性バリエントと良性バリエントを分類する際の、PrimateAIネットワークと他の分類器の正確さを比較する。図30Aの散布図は、DDD症例群vs対照群に対する分類器の各々の性能(y軸)と、保留された霊長類データセットに対する良性予測の正確さ(x軸)とを示す。図30Bは、各分類器に対して示される曲線下面積(AUC)とともに、受信者動作特性(ROC)曲線上に示される、605個の遺伝子内での症例群vs対照群におけるde novoミスセンスバリエントを分離することにおいて異なる分類器を比較する。図30Cは、図28A、図28B、および図28Cにおいて列挙される、PrimateAIネットワークおよび20個の分類器に対する分類の正確さとAUCを示す。示される分類の正確さは、図22Aにおけるエンリッチメントに基づいて予想されるものと同じ数の病原性バリエントと良性バリエントを分類器が予測するような閾値を使用した、真陽性と真陰性の率の平均である。DDD症例群におけるde novoミスセンスバリエントは67%が病原性バリエントであり33%が良性であり、対照群におけるde novoミスセンスバリエントは100%良性であると仮定して、完璧な分類器に対する最大の達成可能なAUCが点線で示されている。

【0405】

図31Aおよび図31Bは、専門家により精選されたClinVarバリエントに対する分類器の性能と、経験的なデータセットに対する性能との相関を示す。図31Aの散布図は、ヒトのみのデータで訓練された、またはヒト+霊長類のデータで訓練された、20個の他の分類器の各々およびPrimateAIネットワークに対する、10000個の保留された霊長類バリエントについての分類の正確さ(x軸)と、ClinVarバリエントについての分類の正確さ(y軸)とを示す。Spearman相関係数rhoおよび関連するP値が示されている。分類器を訓練するために使用されなかったデータへと評価を限定するために、2017年1月から2017年11月の間に追加されたClinVarバリエントのみを使用し、ExAC/gnomADからの一般的なヒトバリエント(>0.1%のアレル頻度)を除外した。示されるClinVar分類の正確さは、ClinVarデータセットにおいて観察されるのと同じ数の病原性バリエントと良性バリエントを分類器が予測するような閾値を使用した、真陽性と真陰性の率の平均である。

【0406】

図31Bの散布図は、ヒトのみのデータで訓練された、またはヒト+霊長類のデータで訓練された、20個の他の分類器およびPrimateAIネットワークの各々に対する、DDD症例群vs対照群の完全なデータセット(x軸)と、ClinVarバリエントに対する分類の正確さ(y軸)を示す。

【0407】

図32は、訓練のための6367個の、妥当性確認のための400個の、および検定のための500個の関連しないタンパク質配列を使用した、Protein DataBankからのアノテートされたサンプルに対する、3状態二次構造予測モデルおよび3状態溶媒接触性予測モデルの性能を示す、補足テーブル14である。配列の相同性が25%未満であるタンパク質のみがProtein DataBankから選択された。3つのクラスは二次構造と溶媒接触性のいずれについても大きく偏ってはいないので、深層学習ネットワークの正確さを性能の尺度として報告する。

【0408】

図33は、予測される二次構造ラベルを使用する深層学習ネットワークとともに、DSSPデータベースからのヒトタンパク質のアノテートされた二次構造ラベルが利用であるときに

10

20

30

40

50



それを使用した、深層学習ネットワークの性能比較を示す補足テーブル15である。

【0409】

図34は、評価した20個の分類器の各々についての、10000個の保留された霊長類バリエーションに対する正確さの値と、DDD症例群vs対照群におけるde novoバリエーションに対するp値とを示す補足テーブル17である。ヒトのデータのみを用いたPrimateAIモデルは、一般的なヒトバリエーション(集団において>0.1%である83500個のバリエーション)のみからなるラベリングされた良性の訓練データセットを使用して訓練された我々の深層学習ネットワークであるが、ヒト+霊長類のデータを用いたPrimateAIモデルは、一般的なヒトバリエーションと霊長類バリエーションの両方を備える、385000個のラベリングされた良性バリエーションの完全なセットについて訓練された我々の深層学習ネットワークである。

10

【0410】

図35は、605個の疾患関連遺伝子に制約された、DDD症例群vs対照群データセットにおけるde novoバリエーションに対する異なる分類器の性能の比較を示す、補足テーブル19である。異なる方法の間で正規化するために、各分類器に対して、DDDセットと対照群セットにおけるエンリッチメントに基づいて予想されるのと同じ数の病原性バリエーションおよび良性バリエーションを分類器が予測するような閾値を特定した。示される分類の正確さは、この閾値における真陽性と真陰性のエラー率の平均である。

【0411】

図49A、図49B、図49C、図49D、および図49Eは、ヒトアレル頻度スペクトラムにわたるミスセンス/同義比を示す。図49Aは、ExAC/gnomADデータベースからの123136人のヒトにおいて観察されるミスセンスバリエーションおよび同義バリエーションが、アレル頻度によって4つのカテゴリへと分割されたことを示す。灰色の影付きの棒は各カテゴリにおける同義バリエーションのカウントを表し、濃緑の棒はミスセンスバリエーションを表す。各棒の高さは各アレル頻度カテゴリにおける同義バリエーションの数に対してスケールされ、ミスセンス/同義カウントおよび比は変異率を調整した後で表示される。図49Bおよび図49Cは、チンパンジーの一般的なバリエーション(図49B)およびチンパンジーのシングルTONバリエーション(図49C)と同一状態(IBS)であるヒトミスセンスバリエーションおよび同義バリエーションに対する、アレル頻度スペクトラムを示す。稀なヒトアレル頻度(<0.1%)と比較して一般的なヒトアレル頻度(>0.1%)を持つチンパンジーミスセンスバリエーションの枯渇率が、付随するカイ二乗(<sup>2</sup>)検定のP値とともに、赤いボックスによって示されている。

20

30

【0412】

図49Dは、ヒト以外の霊長類の種のうちの少なくとも1つにおいて観察されるヒトバリエーションを示す。図49Eは、ClinVarデータベース全体における良性ミスセンスバリエーションと病原性ミスセンスバリエーションのカウント(上の行)を、ExAC/gnomADアレル頻度からサンプリングされる30人のヒトのコホートにおけるClinVarバリエーションのカウント(中間の行)、および霊長類において観察されるバリエーションのカウント(下の行)と比較して示す。矛盾する良性と病原性の判定と、有意性が不確かなものとしてだけアノテートされたバリエーションは、除外された。

【0413】

図50A、図50B、図50C、および図50Dは、他の種と同一状態であるミスセンスバリエーションに対する純化選択を示す。図50Aは、4種の霊長類以外の哺乳類の種(ネズミ、ブタ、ヤギ、およびウシ)において存在するバリエーションと同一状態である、ヒトミスセンスバリエーションおよび同義バリエーションに対するアレル頻度スペクトラムを示す。一般的なヒトアレル頻度(>0.1%)を持つミスセンスバリエーションの枯渇率が、付随するカイ二乗(<sup>2</sup>)検定のP値とともに、赤いボックスによって示されている。

40

【0414】

図50Bは、一般的なヒトアレル頻度(>0.1%)で他の種において観察されるミスセンスバリエーションの枯渇率と、枝長(ヌクレオチドの位置当たりの置換の平均の数)の単位で表される、ヒトからの種の進化的距離とを比較して示す散布図である。各種とヒトとの間の枝長の合計が、種の名前の隣に示されている。親類の個体を含んでいたゴリラを除いて、シング

50

ルトンバリエントおよび一般的なバリエントに対する枯渇値が、バリエント頻度が入手可能であった種に対して示されている。

【0415】

図50Cは、ExAC/gnomADアレル頻度からサンプリングされた30人のヒトのコホートにおける良性ミスセンスバリエントと病原性ミスセンスバリエントのカウント(上の行)を、霊長類において観察されるバリエントのカウント(中間の行)、ならびに、ネズミ、ブタ、ヤギ、およびウシにおいて観察されるバリエントのカウント(下の行)と比較して示す。矛盾する良性と病原性の判定と、有意性が不確かなものとしてだけアノテートされたバリエントは、除外された。

【0416】

図50Dは、一般的なヒトアレル頻度(>0.1%)で近縁の種のペアにおいて観察される固定されたミスセンス置換の枯渇率と、ヒトからの種の進化的距離(平均の枝長の単位で表される)とを比較して示す、散布図である。

【0417】

図51は、純化選択がないときのヒトアレル頻度スペクトラムにわたる予想されるミスセンス:同義比を示す。灰色の影付きの棒は同義バリエントの数を表し、濃緑の棒はミスセンスバリエントの数を表す。点線は同義バリエントによって形成される基準を示す。ミスセンス:同義比は各アレル頻度カテゴリに対して示される。一実装形態によれば、各アレル頻度カテゴリにおける予想されるミスセンスカウントおよび同義カウントは、123136個のエクソンを備えるExAC/gnomADデータセットからイントロンバリエントを取り込み、変異率および遺伝子変換におけるGCバイアスを考慮するバリエントのトリヌクレオチドコンテキストに基づいて、4つのアレル頻度カテゴリの各々に該当することが予想されるバリエントの割合を推定するためにそれらのイントロンバリエントを使用することによって、計算された。

【0418】

図52A、図52B、図52C、および図52Dは、CpGバリエントおよび非CpGバリエントに対するミスセンス:同義比を示す。図52Aおよび図52Bは、ExAC/gnomADエクソンからのすべてのバリエントを使用した、ヒトアレル頻度スペクトラムにわたるCpGバリエントに対するミスセンス:同義比(図52A)および非CpGバリエントに対するミスセンス:同義比(図52B)を示す。図52Cおよび図52Dは、一般的なチンパンジー多型と同一状態であるヒトバリエントのみに制約された、ヒトアレル頻度スペクトラムにわたるCpGバリエントに対するミスセンス:同義比(図52C)および非CpGバリエントに対するミスセンス:同義比(図52D)を示す。

【0419】

図53、図54、および図55は、6種の霊長類と同一状態であるヒトバリエントのミスセンス:同義比を示す。チンパンジー、ボノボ、ゴリラ、オランウータン、アカゲザル、およびマーモセットにおいて存在する変異と同一状態であるExAC/gnomADバリエントに対する、ヒトアレル頻度スペクトラムにわたるミスセンス:同義比のパターンを示す。

【0420】

図56は、調査されるヒトのコホートのサイズを増大させることによって発見される、新しい一般的なミスセンスバリエントの飽和を示すシミュレーションである。シミュレーションにおいて、各サンプルの遺伝子型は、gnomADアレル頻度に従ってサンプリングされた。発見された一般的なgnomADバリエントの割合は、10から100000の各サンプルサイズに対して100回のシミュレーションにわたって平均される。

【0421】

図57は、ゲノムにおける様々な保存プロファイルにわたるPrimateAIの正確さを示す。x軸は、99種の脊椎動物のアラインメントとの、ある配列の周りの51個のAAの百分率のアラインメント可能性を表す。y軸は、10000個の保留された霊長類バリエントの検定データセットについてベンチマークがとられた、保存ピンの各々におけるバリエントに対するPrimateAIの正確さの分類性能を表す。

【0422】

10

20

30

40

50

図58は、一般的なヒトバリエーションおよびヒト以外の霊長類において存在するバリエーションからのラベリングされた良性の訓練データセットへの寄与を示す、補足テーブル5である。

【0423】

図59は、予想されるミスセンス:同義比に対するアレル頻度の影響を示す補足テーブル8である。同義バリエーションとミスセンスバリエーションの予想されるカウントは、変異率および遺伝子変換バイアスを考慮するためにトリヌクレオチドコンテキストを使用して、エクソン境界から少なくとも20~30nt離れたイントロン領域におけるバリエーションのアレル頻度スペクトラムに基づいて計算された。

【0424】

図60は、ClinVar分析を示す補足テーブル9である。一実装形態によれば、ClinVarデータベースの2017年11月のビルドからダウンロードされたバリエーションが、矛盾するアノテーションを伴うミスセンスバリエーションを除外し、有意性が不確かであるバリエーションを削除するようにフィルタリングされ、17775個の良性バリエーションおよび24853個の病原性バリエーションが残った。

【0425】

図61は、一実装形態による、ClinVarにおいて発見された他の種からのミスセンスバリエーションの数を示す補足テーブル10である。バリエーションは、対応するヒトバリエーションと同一状態であることと、同じコーディング結果を保証するためにリーディングフレームの中の他の2つの位置において同一のヌクレオチドを有することとが必要とされた。

【0426】

図62は、元のDDD研究においてはゲノムワイド有意性の閾値にこれまで達していなかった、知的障害における14個の追加の遺伝子候補の発見の一実装形態を示すテーブル1である。

【0427】

図63は、ClinVarにおける病原性バリエーションと良性バリエーションとの間のGranthamスコアの平均の差の一実装形態を示すテーブル2であり、この差は、605個の疾患関連遺伝子内でのDDD症例群vs対照群におけるde novoバリエーションの差の2倍である。

【0428】

[データ生成]

本明細書において使用されるすべての座標は、このセクションで説明される手順を使用して複数配列アラインメントを使用してhg19にマッピングされた他の種におけるバリエーションに対する座標を含めて、ヒトゲノムbuild UCSC hg19/GRCh37を参照する。ヒトとの99種の脊椎動物ゲノムのタンパク質コーディングDNA配列および複数配列アラインメントが、hg19 buildのためのUCSCゲノムブラウザからダウンロードされた(<http://hgdownload.soe.ucsc.edu/goldenPath/hg19/multiz100way/alignments/knownCanonical.exonNuc.fa.gz>)。複数の正規の遺伝子アノテーションを伴う遺伝子については、最長のコーディング転写産物が選択された。

【0429】

世界中の8つの亜集団からの123136人の個人の全エクソンシーケンシング(WES)データを収集した、Exome Aggregation Consortium(ExAC)/genome Aggregation Database(gnomAD) v2.0からヒトエクソン多型データをダウンロードした(<http://gnomad.broadinstitute.org/>)。ExAC VCFファイルにおいてアノテーションされるようなデフォルトの品質制御フィルタを通過しないバリエーション、または正規のコーディング領域の外側にあるバリエーションを除外した。平衡選択による影響を避けるために、霊長類分析のための延長されたMHC領域(chr6:28,477,797~33,448,354)内からのバリエーションも除外した。大型類人猿ゲノムシーケンシングプロジェクトは、24体のチンパンジー、13体のボノボ、27体のゴリラ、および10体のオランウータン(スマトラ亜種からの5体およびボルネオ亜種からの5体を含む、これらを下流分析のために折り畳んだ)に対する、全ゲノムシーケンシングデータおよび遺伝子型を提供する。チンパンジーおよびボノボについての研究は、追加の35体のチンパンジーの

10

20

30

40

50

ゲノム配列を提供する。しかしながら、これらの追加のチンパンジーに対するバリエーションは、大型類人猿ゲノムシーケンシングプロジェクトと同じ方法を使用してコールされなかったため、それらをアレル頻度スペクトラム分析から除外し、深層学習モデルを訓練するためだけに使用した。これらの霊長類多様性研究からの変異はすでに、ヒト基準(hg19)にマッピングされていた。加えて、マーモセットおよびアカゲザルについては、16体のアカゲザルの個体および9体のマーモセットの個体がこれらの種のゲノムの元のシーケンシングにおける変異を評価するために使用されたが、個体レベルの情報は利用可能ではない。

#### 【0430】

大型類人猿ゲノムシーケンシングプロジェクト4は、24体のチンパンジー、13体のボノボ、27体のゴリラ、および10体のオランウータン(スマトラ亜種からの5体およびボルネオ亜種からの5体を含む、これらを下流での分析のために折り畳んだ)に対する、全ゲノムシーケンシングデータおよび遺伝子型を提供する。チンパンジーおよびボノボについての研究は、追加の35体のチンパンジーのゲノム配列を提供する。しかしながら、これらの追加のチンパンジーに対するバリエーションは、大型類人猿ゲノムシーケンシングプロジェクトと同じ方法を使用してコールされなかったため、それらをアレル頻度スペクトラム分析から除外し、深層学習モデルを訓練するためだけに使用した。これらの霊長類多様性研究からの変異はすでに、ヒト基準(hg19)にマッピングされていた。加えて、マーモセットおよびアカゲザルについては、16体のアカゲザルの個体および9体のマーモセットの個体がこれらの種のゲノムの元のシーケンシングにおける変異を評価するために使用されたが、個体レベルの情報は利用可能ではない。

#### 【0431】

他の霊長類および哺乳類と比較するために、アカゲザル、マーモセット、ブタ、ウシ、ヤギ、ネズミ、ニワトリ、およびゼブラフィッシュを含む、他の種のSNPもdbSNPからダウンロードした。dbSNPは追加のオランウータンバリエーションも含んでおり、それを我々は深層学習モデルを訓練するためだけに使用した。それは、個体の遺伝子型情報がアレル頻度スペクトラム分析に対して利用可能ではなかったからである。イヌ、ネコ、またはヒツジなどの他の種については、それらの種に対する限られた数のバリエーションしかdbSNPが提供しないので、廃棄した。

#### 【0432】

バリエーションをヒトにマッピングするために、99種の脊椎動物の多種アラインメントを使用して、ヒトタンパク質コーディング領域へのオーソログ的な1:1のマッピングを確実にした。オーソログ的な多種アラインメントを使用するバリエーションのマッピングが、偽遺伝子またはレトロトランスポジションを受けた配列によって引き起こされるアーティファクトを取り除くために必須であった。このアーティファクトは、多数対1のマッピングを可能にするliftOverなどのツールを使用して種間で直接SNPをマッピングするときには発生する。dbSNPにおける種のゲノムビルドが99種の脊椎動物の複数配列アラインメントにおける種のゲノムビルドと一致しなかった場合、複数配列アラインメントにおいて使用されるゲノムビルドへとバリエーションを更新するためにliftOverを使用した。バリエーションが基準/代替方向において発生した場合、バリエーションを同一状態であるものとして受け入れた。たとえば、ヒト基準がGであり代替アレルがAであった場合、これは、基準がAであり代替アレルがGであった別の種におけるバリエーションと同一状態であると見なされた。バリエーションがヒトと他の種の両方において同じ予測されるタンパク質コーディング結果を有することを確実にするために、ミスセンスバリエーションと同義バリエーションの両方に対して、コドンの中の他の2つのヌクレオチドが種間で同一であることを要求した。分析に含まれる各種からの多型が補足データファイル1において列挙され、詳細な尺度が補足テーブル1に示されている。

#### 【0433】

各dbSNP提出者バッチからのバリエーションが、高品質でありヒトに正しくアラインメントされることを確実にするために、各バッチに対するミスセンス:同義比を計算し、これが2.2:1という予想される比より小さかったことと、大半の種、特に非常に大きい有効個体数

10

20

30

40

50

を有することが予想されるゼブラフィッシュおよびネズミが、1:1より小さい比を有していたことを確認した。さらなる分析から、異常に高いミスセンス:同義比を有していたウシから、SNPの2つのバッチを除外した(比が1.391であるsnpBatch\_1000\_BULL\_GENOMES\_1059190.gzおよび比が2.568であるsnpBatch\_COFACTOR\_GENOMICS\_1059634.gz)。残りのウシのバッチに対する平均のミスセンス:同義比は0.8:1であった。

#### 【0434】

[ミスセンス:同義比、変異率、遺伝的浮動、およびGCバイアス(GC-Biased)遺伝子変換に対するアレル頻度の影響の訂正]

純化選択の活動に加えて、高いアレル頻度でのヒトミスセンスバリエーションの観察される枯渇率は、自然選択に関連しない要因によっても影響を受け得る。集団において特定のアレル頻度で現れる自然変異の確率は、変異率、遺伝子変換、および遺伝的浮動の関数であり、これらの要因は、選択圧がなくてもアレル頻度スペクトラムにわたってミスセンス:同義比にバイアスをもたらす可能性がある。

#### 【0435】

タンパク質コーディング選択がないときの各アレル頻度カテゴリにおける予想されるミスセンス:同義比を計算するために、各エクソンの31~50bp上流および21~50bp下流のイントロン領域内のバリエーションを選択した。これらの領域は、延長されたスプライスマチーフの影響を避けるのに十分遠くなるように選ばれた。これらの領域は、ExAC/gnomADエクソンに対するエクソン捕捉配列の端に近いので、バリエーションの公平な確認を確実にするために、あらゆるchrX領域を除去し、平均リード深さが30未満である領域を取り除いた。各バリエーションおよびそのすぐ上流および下流のヌクレオチドは、64個のトリヌクレオチドコンテキストのうちの一つの中にある。中間のヌクレオチドを3つの他の塩基へと変異させる場合、全体で $64 \times 3 = 192$ 個のトリヌクレオチド構成が可能である。トリヌクレオチド構成およびそれらの逆相補配列は等価であるので、実質的に96個のトリヌクレオチドコンテキストがある。トリヌクレオチドコンテキストは変異率に対して非常に強い影響があり、GCバイアス遺伝子変換に対する影響がより小さいことが観察され、これにより、これらの変数をモデル化するためにトリヌクレオチドコンテキストが有効になる。

#### 【0436】

これらのイントロン領域内で、アレル頻度の4つのカテゴリ(シングルトン、シングルトンより多い~0.01%、0.01~0.1%、>0.1%)および192個のトリヌクレオチドコンテキストに基づいて、126136個のExAC/gnomADエクソンから各バリエーションを取り込み、それらを $4 \times 192$ 個のカテゴリへと分離した。そのトリヌクレオチドコンテキスト(イントロン配列の中の各ヌクレオチドを3つの異なる方式で置換することによって得られる)で潜在的なバリエーションの総数を割ることによって、 $4 \times 192$ 個のカテゴリ(アレル頻度 $\times$ トリヌクレオチドコンテキスト)の各々において観察されるバリエーションの数を正規化した。192個のトリヌクレオチドコンテキストの各々に対して、タンパク質コーディング選択がないとき、4つのアレル頻度カテゴリの各々に該当するバリエーションの予想される割合をこうして得た。これは、トリヌクレオチドコンテキストの差による、変異率、GCバイアス遺伝子変換、および遺伝的浮動の影響を暗黙的にモデル化する(補足テーブル7)。

#### 【0437】

各アレル頻度カテゴリにおける予想されるミスセンス:同義比を得るために、一塩基置換によって入手可能なヒトゲノムにおける潜在的な同義変異およびミスセンス変異の総数をカウントし、それらの各々を192個のトリヌクレオチドコンテキストのうちの一つに割り当てた。各コンテキストに対して、4つのアレル頻度カテゴリの各々に該当することが予想されるバリエーションの数を計算するために、 $4 \times 192$ 個のテーブルを使用した。最後に、192個のトリヌクレオチドコンテキストにわたる同義バリエーションおよびミスセンスバリエーションの数を合計して、4つのアレル頻度カテゴリの各々における同義バリエーションとミスセンスバリエーションの合計の予想される数を得た(図51および補足テーブル8(図59))。

#### 【0438】

予想されるミスセンス:同義比が2.46:1であったシングルトンバリエーションを除き、予想

10

20

30

40

50

されるミスセンス:同義比は、アレル頻度スペクトラムにわたってほぼ一定であり、自然選択がない場合にde novoバリエーションについて予想される2.23:1という比に近かった。このことは、タンパク質コーディング選択圧力とは独立の要因(変異率、遺伝子変換、遺伝的浮動)の活動により、ExAC/gnomADにおけるシングルtonアレル頻度カテゴリのバリエーションが、デフォルトでde novo変異より約10%高いミスセンス:同義比を有することが予想されることを示す。これを訂正するために、アレル頻度分析において、シングルtonに対するミスセンス:同義比を10%だけ下に調整した(図49A、図49B、図49C、図49D、および図49E、ならびに図50A、図50B、図50C、および図50D)。この小さい調整は、霊長類および他の哺乳類において存在する一般的なヒトバリエーションに対する推定されるミスセンス枯渇率をおよそ3.8%低下させた(図49A、図49B、図49C、図49D、および図49E、ならびに図50A、図50B、図50C、および図50Dに示される)。シングルtonバリエーションに対するより高いミスセンス:同義比は、トランスポージョン変異(これはミスセンス変化を作り出す可能性がより高い)より高い変異率が原因で、トランジション変異(これは同義変化を作り出す可能性がより高い)がより高いアレル頻度を有することが原因である。

#### 【0439】

その上、このことは、2.23:1というde novo変異について予想される比を超える、ExAC/gnomADにおけるシングルtonバリエーションに対する2.33:1という観察されるミスセンス:同義比を説明する。ミスセンス:同義比に対するアレル頻度スペクトラムの影響を考慮した後で、これは実際には予想と比較して5.3%のシングルtonバリエーションの枯渇を反映しており、これはおそらく、de novo優性遺伝モードを持つ病原性ミスセンス変異に対する選択によるものである。実際に、機能喪失の確率が高い(pLI>0.9)ハプロ不全遺伝子だけを考慮するとき、ExAC/gnomADシングルtonバリエーションに対するミスセンス:同義比は2.04:1であり、ハプロ不全遺伝子内で約17%前後の枯渇率を示す。この結果は、ある程度の不完全な浸透を仮定するとき、ミスセンス変異の20%が機能喪失変異と等価であるという以前の推定と合致する。

#### 【0440】

変異率の大きな違いによる、ヒトアレル頻度スペクトラムにわたるCpGバリエーションおよび非CpGバリエーションに対するミスセンス:同義比も具体的に調査した(図52A、図52B、図52C、および図52D)。CpG変異と非CpG変異の両方に対して、一般的なチンパンジー多型と同一状態であるヒトバリエーションは、アレル頻度スペクトラムにわたってほぼ一定のミスセンス:同義比を有することを確認した。

#### 【0441】

[他の種における多型と同一状態であるヒトミスセンスバリエーションの枯渇率]

他の種からのバリエーションがヒトにおいて一般的なアレル頻度(>0.1%)で耐えられるかどうかを評価するために、他の種における変異と同一状態であったヒトバリエーションを特定した。バリエーションの各々に対して、ヒト集団におけるアレル頻度(シングルton、シングルtonより多い~0.01%、0.01%~0.1%、>0.1%)に基づいて、それらを4つのカテゴリのうち1つに割り当て、稀なバリエーション(<0.1%)と一般的なバリエーション(>0.1%)との間でのミスセンス:同義比(MSR)の低下を推定した。一般的なヒトアレル頻度(>0.1%)における同一状態のミスセンスバリエーションの枯渇率は、ヒトにおいて一般的なアレル頻度では自然選択により除去されるのに十分有害な、他の種からのバリエーションの割合を示す。

#### 【0442】

##### 【数42】

$$\% \text{depletion} = \frac{\text{MSR}_{\text{rare}} - \text{MSR}_{\text{comm}}}{\text{MSR}_{\text{rare}}}$$

#### 【0443】

ミスセンス:同義比および枯渇の百分率は、種毎に計算され、図50Bおよび補足テーブル

2に示されている。加えて、一般的なチンパンジーバリエーション(図49A)、チンパンジーシングルバリエーション(図49C)、および哺乳類バリエーション(図50A)に対して、稀なバリエーションと一般的なバリエーションとの間のミスセンス:同義比の差が有意であったかどうかを検定するために、2×2の分割表上で相同性のカイ二乗(2)検定を実行した。

#### 【0444】

シーケンシングは大型類人猿多様性プロジェクトからの限られた数の個体に対してのみ実行されたので、一般のチンパンジー集団において稀(<0.1%)または一般的(>0.1%)であったサンプリングされたバリエーションの割合を推定するために、ExAC/gnomADからのヒトアレレル頻度スペクトラムを使用した。ExAC/gnomADアレレル頻度に基づく24体の個体のコホートをサンプリングし、このコホートにおいて一度観察された、または一度より多く観察されたミスセンスバリエーションを特定した。一度より多く観察されたバリエーションは、一般の集団において一般的(>0.1%)である確率が99.8%であったのに対し、コホートにおいて一度だけ観察されたバリエーションは、一般の集団において一般的である確率が69%であった。図49Bおよび図49Cにおいて、チンパンジーシングルバリエーションの一部が稀な有害な変異であることの結果として、ヒトにおいて高いアレレル頻度でシングルチンパンジーバリエーションの枯渇が観察されるが、一般的なチンパンジーバリエーションについてはそれが観察されないということが示される。24体の個体のコホートにおいて観察されるチンパンジーバリエーションの概ね半分は一度だけ観察され、概ね半分は一度より多く観察された。

#### 【0445】

より遠縁の哺乳類におけるミスセンスバリエーションに対する観察される枯渇率は、よく保存されている、したがってより正確にアラインメントされている遺伝子の混乱をもたらす影響によるものではなかったことを確認するために、ヒトと比較して11種の霊長類および50種の哺乳類の複数配列アラインメントにおける50%を超える平均ヌクレオチド相同性を持つ遺伝子のみを制約して、上記の分析を繰り返した(補足テーブル3参照)。これは、結果に実質的な影響を与えることなく、分析から約7%のヒトタンパク質コーディング遺伝子を除去した。

#### 【0446】

[ 霊長類、哺乳類、および遠縁の脊椎動物の間で固定された置換 ]

バリエーションデータについての問題、または家畜化によるアーティファクト(dbSNPから選択された種の大半は家畜化されているので)により、dbSNP変異を使用した我々の結果が影響を受けなかったことを確実にするために、種内多型の代わりに近縁の種のペアからの固定された置換を使用した分析も繰り返した。枝長で測定される進化系統距離(場所当たりのヌクレオチド置換の平均の数)とともに、UCSCゲノムブラウザから100種の脊椎動物の種の進化系統樹をダウンロードした(<http://hgdownload.soe.ucsc.edu/goldenPath/hg19/multiz100way/hg19.100way.commonNames.nh>)。さらなる分析のために、近縁の種のペア(枝長<0.25)を選択した。近縁の種のペア間の固定された置換を特定するために、ヒトとの99種の脊椎動物ゲノムの複数配列アラインメントのための、ならびにヒトとの19種の哺乳類(16種の霊長類)ゲノムのアラインメントのための、コーディング領域をUCSCゲノムブラウザからダウンロードした。追加の19種の哺乳類の複数種アラインメントは、ボノボなどの霊長類の種の一部が99種の脊椎動物アラインメントにおいて存在しなかったので必要であった(<http://hgdownload.soe.ucsc.edu/goldenPath/hg38/multiz20way/alignments/knownCanonical.exonNuc.fasta.gz>)。全体で、図50Dおよび補足テーブル4に列挙されるように、5つの霊長類ペアを含む、近縁の種の15個のペアを得た。

#### 【0447】

正規のコーディング領域内でのヒトとの19種の哺乳類ゲノムまたは99種の脊椎動物ゲノムの複数配列アラインメントを取り込み、補足データファイル2において列挙される、脊椎動物の各々の選択されたペア間でのヌクレオチド置換を得た。コドンの中の他の2つのヌクレオチドがヒトと他の種との間で変わらなかったことを条件とし、かつ基準方向と代替方向のいずれかのバリエーションを受け入れて、これらの置換がヒトゲノムにマッピングされた。関連する種のペアからの固定された置換と同一状態であったヒトバリエーションを使用し

て、稀な(<0.1%)アレル頻度カテゴリのパリアントと一般的な(>0.1%)アレル頻度カテゴリのパリアントに対するミスセンス:同義比を計算し、補足テーブル4において示されるように、負の選択のもとで固定された置換の割合を得た。

#### 【0448】

[ヒト、霊長類、哺乳類、および他の脊椎動物に対する多型データのClinVar分析]

他の種と同一状態であるパリアントの臨床上的影響を調査するために、ClinVarデータベース(2017年11月2日に発表されたftp://ftp.ncbi.nlm.nih.gov/pub/clinvar/clinvar\_20171029.vcf.gz)12のリリースパリアントサマリ(release variant summary)をダウンロードした。このデータベースは、hg19ゲノムビルド上の324698個のパリアントを含み、そのうち122884個がタンパク質コーディング遺伝子の我々のリストにマッピングするミスセンス塩基パリアントであった(補足テーブル9)。ClinVarデータベースの中のパリアントの大半はミスセンスの結果をもたらさず、除外された。次に、矛盾する病原性の解釈を持つパリアントをフィルタリングし、良性、良性である可能性が高い、病原性、および病原性である可能性が高いアノテーションを伴うパリアントのみを残した。良性のアノテーションおよび良性である可能性が高いというアノテーションを持つパリアントを単一のカテゴリへと統合し、病原性のアノテーションまたは病原性である可能性が高いというアノテーションを持つパリアントも統合した。補足テーブル9に示されるフィルタリングステップの後で、全体で病原性カテゴリの中の24853個のパリアントおよび良性カテゴリの中の17775個のパリアントがあり、残りは有意性が知られていないまたは矛盾するアノテーションを伴うパリアントであるので、除外された。

10

20

#### 【0449】

ヒト集団におけるClinVarミスセンスパリアントに対する基準を得るために、ExAC/gnomADアレル頻度からサンプリングされた30人の個人の cohorts においてClinVarミスセンスパリアントを調査した。この cohort サイズは、霊長類多様性プロジェクト研究においてシーケンシングされた個人の数を概ね反映するように選ばれた。100個のそのようなシミュレーションからの、30人のヒトの cohorts における病原性パリアントと良性パリアントの平均の数を報告する(図49E)。専門家は、ClinVarにおいて一般的なヒトパリアントを良性の結果で系統的にアノテートしてきたので、この精選のバイアスを避けるためにアレル頻度が1%より高いパリアントを除外した。

30

#### 【0450】

霊長類、哺乳類、および他の脊椎動物における変異と同一状態であったClinVarパリアントを分析した。各種に対する良性パリアントおよび病原性パリアントの数が補足テーブル10に示されている。ヒト、霊長類、およびより遠縁の哺乳類において存在したClinVarパリアントの数の概要が、良性パリアントと病原性パリアントの比の差についての相同性のカイ二乗(2)検定からの結果とともに、図49Eおよび図50Bにおいて示されている。

#### 【0451】

[モデル訓練のための良性パリアントの生成]

ヒト集団において一般的なパリアントは、創始者効果または平衡選択の稀な事例を除いて大部分が中立的であり、これにより、それらのパリアントは、人の解釈によるバイアスの影響を受けていない機械学習のための良性訓練データセットとして適切なものになる。フィルタを通過しなかったパリアントを除いて、ExAC/gnomADデータベース(リリースv2.0)からの123136個のエクソンからアレル頻度データを使用し、正規のタンパク質コーディング転写産物内で全体の集団アレル頻度が0.1%以上である83546個のミスセンスパリアントが残った。

40

#### 【0452】

霊長類において存在するパリアントは大部分がヒトにおいて良性であることを示す我々の以前の結果に基づいて、一般的なヒトパリアント(>0.1%のアレル頻度)、大型類人猿多様性プロジェクトおよび追加の霊長類シーケンシングからのチンパンジー、ボノボ、ゴリラ、およびオランウータンからのパリアント、ならびに、dbSNPからのアカゲザル、オランウータン、およびマーモセットのパリアントを備える、機械学習のための良性訓練デー

50



タセットを作成した。一実装形態によれば、全体で、良性訓練セットに301690個の固有の霊長類バリエーションが追加された。各源が寄与した良性訓練バリエーションの数が補足テーブル5に示されている。

#### 【0453】

注意すべき点は、大半の霊長類バリエーションはそれらのそれぞれの集団において一般的であり、少数派が稀なバリエーションであるということである。ヒト以外の霊長類の種は、シーケンシングされた個体の数が限られていたので、確認されたバリエーションのセットは全般に、一般的な変異を表すことが予想される。実際に、霊長類の種の各々からのバリエーションに対するミスセンス:同義比は、*de novo*変異に対する予想される2.23:1の比の半分未満であることが見出され、これらの大半が選択のふるいにすでにかけてきた一般的なバリエーションであることを示している。その上、チンパンジーのコホートに対して、確認されたバリエーションの約84%が、それらのそれぞれの集団において一般的なアレル頻度(>0.1%)で存在することが推定された。新しく見つかったミスセンス変異の約50%が一般的なヒトアレル頻度(>0.1%)において純化選択によってフィルタリングされるので(図49A)、この数字は、観察された霊長類変異と同一状態であるヒトミスセンスバリエーションの8.8%という観察された枯渇率を説明する、約16%の稀なバリエーションと一致している(図49D)。

#### 【0454】

ヒトミスセンス変異の約20%が機能喪失と等価であるという推定を適用すると、霊長類バリエーションは3.2%の完全に病原性の変異と、91.2%の良性の変異(>0.1%のアレル頻度に耐える)と、完全には遺伝子の機能を失わず一般的なアレル頻度(>0.1%)で除去されるほど十分に有害ではない5.6%の中間的な変異とを含むことが予想される。この訓練データセットは不完全であることが知られているが、深層学習ネットワークの分類の正確さは、一般的なヒトバリエーションと霊長類バリエーションの両方を備える良性訓練データセット上で訓練されたとき、一般的なヒトバリエーションのみの場合と比較してはるかに良好であった。したがって、現在の分類の正確さでは、利用可能な訓練データの量はより強い制限であるように見える。より多数の個体が各霊長類の種においてシーケンシングされるにつれて、より高い割合の一般的な霊長類バリエーションを含む訓練データセットを準備することが可能になり、訓練データセットにおける病原性バリエーションからの汚染が減り、分類性能がさらに改善する。

#### 【0455】

[ 良性訓練データセットを補足するためのラベリングされていないバリエーションの生成 ]

すべての潜在的なミスセンスバリエーションが、正規のコーディング領域の各塩基場所から、その場所におけるヌクレオチドを他の3つのヌクレオチドで置換することによって生成された。ExAC/gnomADからの123136個のエクソンにおいて観察されたバリエーションと、開始コドンまたは終止コドンにおけるバリエーションを除外した。全体で、68,258,623個のラベリングされていないバリエーションが生成された。ラベリングされていないバリエーションの各々を、96個の異なるトリヌクレオチドコンテクストカテゴリのうちの一つに割り当てた。トリヌクレオチドコンテクストによって良性データセットの中のバリエーションと一致する、このラベリングされていないデータセットからバリエーションをサンプリングし、良性の訓練例とラベリングされていない訓練例を区別するように分類器を訓練することによって、半教師ありの手法を使用して深層学習ネットワークを訓練した。

#### 【0456】

[ ラベリングされていないバリエーションの追加のフィルタリング ]

良性バリエーションおよびラベリングされていないバリエーションの例をランキングアミノ酸配列とともに提示することによって、深層学習ネットワークは、変異に対して高度に耐性のないタンパク質の領域を学習する。しかしながら、タンパク質配列の領域に一般的なバリエーションがないことは、強い純化選択によるものであることがあり、または、バリエーションが領域においてコールされるのを妨げる技術的なアーティファクトによるものであることがある。後者を訂正するために、ExAC/gnomADデータセットが1より小さい平均カバレッジを有していた領域から、良性データセットとラベリングされていないデータセットの両方

からのバリエーションを除去した。同様に、ラベリングされていないバリエーションを訓練の間に良性データセットの中の霊長類バリエーションと照合するとき、霊長類が複数配列アラインメントにおいてヒトとのオーソロガスなアラインメント可能な配列を有しなかった領域から、ラベリングされていないバリエーションを除外した。

#### 【0457】

[ 妥当性確認および検定のための保留された霊長類バリエーション、ならびに影響を受けている個人および影響を受けていない個人からのde novoバリエーション ]

深層学習ネットワークの妥当性確認および検定のために、妥当性確認および検定のために10000個の霊長類バリエーションの2つのセットをランダムにサンプリングし、これらについては訓練を保留した。霊長類バリエーションの残りは、一般的なヒトバリエーション(>0.1%のアレル頻度)とともに、深層学習ネットワークを訓練するための良性データセットとして使用された。加えて、妥当性確認セットおよび検定セットのために、保留された霊長類バリエーションと照合された10000個のラベリングされていないバリエーションの2つのセットもサンプリングした。

#### 【0458】

2つのセットの中のバリエーションを区別するためのネットワークの能力を測ることによって、訓練の過程の間に深層学習ネットワークの性能を監視するために、妥当性確認セットの中の10000個の保留された霊長類バリエーションを使用し、10000個のラベリングされていないバリエーションを照合した。これにより、ネットワークの性能が飽和すると、訓練の停止点を決定して、過剰適応を避けることが可能になった。

#### 【0459】

深層学習ネットワークならびに他の20個の分類器のベンチマークをとるために、検定データセットの中の10000個の保留された霊長類バリエーションを使用した。異なる分類器は大きく変化するスコア分布を有していたので、これらのラベリングされていないバリエーションを使用して、各分類器に対する50パーセントイル閾値を特定した。方法間の公平な比較を確実にするために、その分類器に対して50パーセントイル閾値で良性であるものと分類された、10000個の保留された霊長類バリエーション検定セットの中のバリエーションの割合について、各分類器のベンチマークをとった。

#### 【0460】

神経発達障害を持つ影響を受けている個人におけるde novoバリエーションと、健康な対照群におけるde novoバリエーションとを使用して、臨床上の環境において深層学習ネットワークの性能を評価するために、Deciphering Developmental Disorders(DDD)研究からのde novoバリエーションと、Simons Simplex Collection(SSC)自閉症研究における健康な兄弟の対照群からのde novoバリエーションとをダウンロードした。DDD研究はde novoバリエーションに対する信頼性レベルを提供し、バリエーションコーリングエラーによる潜在的な偽陽性として、閾値が0.1未満であるDDDデータセットからのバリエーションを除外した。全体で、DDDの影響を受けている個人からの3512個のミスセンスde novoバリエーションおよび健康な対照群からの1208個のミスセンスde novoバリエーションがあった。

#### 【0461】

疾患遺伝子候補の顔ぶれの中で、有意性が不確かな良性バリエーションおよび病原性バリエーションを区別するという現実世界の臨床上のシナリオをより良くモデル化するために、タンパク質切断変異だけから計算された(補足テーブル18)、DDD研究において疾患と関連付けられた( $p < 0.05$ )605個の遺伝子内のde novoバリエーションのみに分析を限定した。遺伝子固有の変異率および考慮される染色体の数を仮定して、de novo変異の予想される数のヌル仮説のもとで統計的有意性を計算することによって、タンパク質切断de novo変異の遺伝子固有のエンリッチメントを評価した。名目のP値が0.05未満である605個の遺伝子を選択した。605個の遺伝子内の同義de novo変異およびミスセンスde novo変異の余剰(図22A)を、観察されるde novo変異vs予想されるde novo変異のカウントの比として、ならびに、観察されるde novo変異から予想されるde novo変異を引いた差分として計算した。これらの605個の遺伝子内で、DDDの影響を受けている個人から380個のde novoミスセンス変異を観察

10

20

30

40

50

した(図22A)。我々自身の分類器を含む分類器の各々に対して、小さい割合のバリエーションは予測を有せず、それは全般に、それらのバリエーションが分類器によって使用される同じ転写産物モデルにマッピングしなかったからである。したがって、我々の深層学習ネットワークでは、DDDの影響を受けている個人からの362個のde novoミスセンス変異および健康な対照群からの65個のde novoミスセンス変異を使用して、図22A、図22B、図22C、図22D、および図22Eにおいて下流分析を実行した。

#### 【0462】

[シーケンシングされた霊長類集団の数の増大に伴うすべての潜在的なヒトミスセンス変異の飽和]

504種の現存する霊長類の種において存在する一般的なバリエーションによる、すべての700万個の潜在的なヒトミスセンス変異の予想される飽和を調査した。各霊長類の種に対して、ヒトにおいて観察される一般的なミスセンスバリエーションの数(アレル頻度が0.1%より高い約83500個のミスセンスバリエーション)を4回シミュレートした。それは、他の霊長類の種と比べてヒトの個体当たりのバリエーションの数が概ね半分であるように見え、ヒトミスセンスバリエーションの約50%が0.1%を超えるアレル頻度において純化選択により除去されているからである(図49A)。96個のトリヌクレオチドコンテキストにおける一般的なヒトミスセンスバリエーションの観察される分布に基づいて、シミュレートされたバリエーションを割り当てた。たとえば、一般的なヒトミスセンスバリエーションの2%が、CCC>CTGのトリヌクレオチドコンテキストからのものであった場合、シミュレートされるバリエーションの2%がランダムにサンプリングされたCCG>CTG変異であったことを要求した。これは、トリヌクレオチドコンテキストを使用して、変異率、遺伝的浮動、および遺伝子変換バイアスの影響を考慮する効果を有する。

#### 【0463】

図23Dの曲線は、各霊長類の種におけるすべての一般的なバリエーション(>0.1%のアレル頻度)を確認していることを仮定して、504種の霊長類の種のいずれかに存在する一般的なバリエーションによる約7000万個の潜在的なヒトミスセンス変異の累積的な飽和を示す。図49Aから、ヒトミスセンス変異の概ね約50%が、ヒトと他の霊長類の両方において、一般的なアレル頻度(>0.1%)に達するのを妨げるのに十分有害であるので、図23Dの曲線は、無害なヒトミスセンス変異の割合が、霊長類の種の数が増えるにつれて一般的な霊長類変異により飽和することを表す。504種の霊長類の種を用いると、無害なヒトミスセンス変異の大半が飽和し、無害なCpG変異は変異率がより高いことが原因ではるかに少数の種で飽和する。

#### 【0464】

調査されたヒトのコホートのサイズを大きくすることで発見された一般的なヒトミスセンスバリエーション(>0.1%のアレル頻度)の割合をモデル化するために(図36)、gnomADアレル頻度に従って遺伝子型をサンプリングした。発見された一般的なgnomADミスセンスバリエーションの割合は、10から100000までの各サンプルサイズに対する100回のシミュレーションにわたって平均された。

#### 【0465】

[二次構造および溶媒接触性の予測]

病原性予測のための深層学習ネットワークは、二次構造および溶媒接触性予測ネットワークのための19個の畳み込み層と、二次構造および溶媒接触性ネットワークの結果を入力として取り込む主病原性予測ネットワークのための17個の畳み込み層とを含む、全体で36個の畳み込み層を含む。大半のヒトタンパク質の結晶構造は知られていないので、ネットワークが一次配列からタンパク質構造を学習することを可能にするために2つのモデルを訓練した。両方のモデルが、図6に示される同じネットワークアーキテクチャおよび入力を使用した。二次構造および溶媒接触性ネットワークへの入力は、99種の他の脊椎動物とのヒトの複数配列アラインメントからの保存情報を符号化する、長さ51×20個のアミノ酸の位置特定の頻度行列である。

#### 【0466】

10

20

30

40

50

二次構造ネットワークは、ヘリックス(H)、シート(B)、およびコイル(C)という3状態の二次構造を予測するように訓練される。溶媒接触性ネットワークは、埋もれている(B)、中間(I)、および露出している(E)という3状態の溶媒接触性を予測するように訓練される。両方のネットワークが一次配列のみを入力として取り込み、Protein DataBankにおける既知の結晶構造からのラベルを使用して訓練された。モデルは各アミノ酸残基に対して1つの状態を予測する。

#### 【0467】

[二次構造および溶媒接触性の予測のためのデータ準備]

モデルを訓練するために、Protein Databankからの関連しない結晶構造を使用した。25%を超える配列相動性を持つアミノ酸配列が除去された。全体で、6367個のタンパク質配列が訓練のために使用され、400個が妥当性確認のために使用され、500個が検定のために使用された(補足テーブル13)。アミノ酸配列および二次構造と溶媒接触性ラベルを含む、訓練のために使用されたデータは、RaptorXウェブサイト:<http://raptorx.uchicago.edu/download/>から入手可能である。

10

#### 【0468】

大半の解かれている結晶構造はヒト以外のタンパク質のものであるので、二次構造および溶媒モデルを事前訓練するために、RaptorXスイート(PSI-BLASTに基づく)を使用して関連する配列を取得した。それは、ヒトベースの複数配列アラインメントが一般に入手可能ではなかったからである。RaptorXからCNFsearch1.66\_releaseツールを使用してタンパク質に対する複数配列アラインメントを生成し、99個の最も近いアラインメントから各場所におけるアミノ酸をカウントして位置特定の頻度行列を形成した。たとえば、1u71A.fastaタンパク質に対する複数配列アラインメントを読み出すためのRaptorXを使用する具体的なコマンドは次の通りであった。

20

```
% ./buildFeature-i 1u71A.fasta-c 10-o ./TGT/1u71A.tgt
% ./CNFsearch-a 30-q 1u71A
```

#### 【0469】

データセットの中の各アミノ酸の場所に対して、51個のランキングアミノ酸に対応する位置特定の頻度行列からウィンドウを取り込み、これを使用して長さ51のアミノ酸配列の中心にあるアミノ酸に対する二次構造または溶媒接触性のいずれかのためのラベルを予測した。二次構造および相対的な溶媒接触性のためのラベルは、DSSPソフトウェアを使用してタンパク質の既知の3D結晶構造から直接取得され、一次配列からの予測を必要としなかった。二次構造および溶媒接触性ネットワークを病原性予測ネットワークの一部として組み込むために、ヒトベースの99種の脊椎動物の複数配列アラインメントから位置特定の頻度行列を計算した。これらの2つの方法から生成された保存行列は一般に類似しているが、パラメータ重みの精密な調整を可能にするために、病原性予測のための訓練の間に二次構造および溶媒接触性モデルを通じた逆伝播を可能にした。

30

#### 【0470】

[モデルアーキテクチャおよび訓練]

タンパク質の二次構造および相対的な溶媒接触性を予測するように、2つの別々の深層畳み込みニューラルネットワークモデルを訓練した。2つのモデルは、同一のアーキテクチャおよび入力データを有するが、予測状態については異なる。最高の性能に向けてモデルを最適化するために、詳細なハイパーパラメータ探索を行った。病原性予測のための我々の深層学習ネットワークと、二次構造および溶媒接触性を予測するための深層学習ネットワークの両方が、画像分類における成功により広く採用されている残差ブロックのアーキテクチャを採用した。残差ブロックは、より前の層からの情報が残差ブロックをスキップすることを可能にするスキップ接続が散在する、反復する畳み込みのユニットを備える。各残差ブロックにおいて、入力層がまずバッチ正規化され、正規化線形ユニット(ReLU)を使用する活性化層がそれに続く。活性化は次いで1D畳み込み層を通される。1D畳み込み層からのこの中間の出力は、再びバッチ正規化およびReLU活性化され、別の1D畳み込み層がそれに続く。第2の1D畳み込みの終わりに、その出力を元の入力と合計して残差ブロッ

40

50

クにし、このことが、元の入力情報が残差ブロックをバイパスすることを可能にすることによってスキップ接続として活動する。著者により深層残差学習ネットワークと名付けられるそのようなアーキテクチャでは、入力はその元の状態で保存され、残差接続にはモデルからの非線形の活性化がない状態に保たれ、より深いネットワークの効果的な訓練が可能になる。詳細なアーキテクチャは、図6および補足テーブル11(図7Aおよび図7B)および図12(図8Aおよび図8B)において提供される。

#### 【0471】

残差ブロックに続いて、ソフトマックス層が各アミノ酸に対する3状態の確率を計算し、それらの中で最大のソフトマックス確率がアミノ酸の状態を決定する。モデルは、ADAM最適化器を使用して、タンパク質配列全体に対する累積カテゴリークロスエントロピー損失関数(accumulated categorical cross entropy loss function)を用いて訓練される。ネットワークが二次構造および溶媒接触性について事前訓練されると、病原性予測ネットワークのための入力としてネットワークの出力を直接取り込む代わりに、ソフトマックス層の前の層を取り込み、それによって、より多くの情報が病原性予測ネットワークを通る。

#### 【0472】

3状態二次構造予測モデルについて達成される最高の検定の正確さは79.86%(補足テーブル14)であり、DeepCNF model30により予測される最新の正確さと同様である。3状態溶媒接触性予測モデルに対する最高の検定の正確さは60.31%(補足テーブル14)であり、同様の訓練データセット上でRaptorXによって予測される現在の最高の正確さと同様である。結晶構造を有していた約4000個のヒトタンパク質に対するDSSPでアノテートされた構造ラベルを使用するとき、予測された構造ラベルのみを使用する標準的なPrimateAIモデルを使用するときの、ニューラルネットワークの予測の比較も行った。DSSPでアノテートされたラベルを使用するとき、病原性予測の正確さにおけるさらなる改善は得られなかった(補足テーブル15)。

#### 【0473】

[病原性予測のための深層学習モデルの入力特徴量]

病原性予測ネットワークのための訓練データセットは、フィルタリングの後で、385236個の良性とラベリングされたバリエーションと、68258623個のラベリングされていないバリエーションを含む。各バリエーションに対して、以下の入力特徴量を生成した。各バリエーションの第1の入力特徴量は、バリエーションの配列コンテキストを深層学習モデルに提供するための、長さ51のランキングアミノ酸配列、すなわち、hg19の基準配列から得られたバリエーションの各側への25個のアミノ酸である。全体で、このランキング基準配列は長さが51個のアミノ酸である。経験的な観察結果を通じて、タンパク質配列のアミノ酸表現が、ヌクレオチドを使用してタンパク質コーディング配列を表現することより効果的であったことを発見した。

#### 【0474】

第2の特徴量は、バリエーションによって中心の場所において置換された代替アミノ酸を伴う、長さ51のランキングヒトアミノ酸配列である。代替ランキング配列は、配列の中央の場所が基準アミノ酸の代わりに代替アミノ酸を含むことを除き、第1の特徴量における基準ランキング配列と同じである。基準ヒトアミノ酸配列と代替ヒトアミノ酸配列の両方が、長さ51×20のワンホット符号化されたベクトルへと変換され、各アミノ酸は、値0を伴う19個のアミノ酸および値1を伴う単一のアミノ酸のベクトルによって表される。

#### 【0475】

11種の霊長類のための1つの位置特定の頻度行列(PFM)、霊長類を除く50種の哺乳類のための1つのPFM、および霊長類と哺乳類を除く38種の脊椎動物のための1つのPFMを含む、3つのPFMが、バリエーションに対する99種の脊椎動物の複数配列アラインメントから生成される。各PFMはL×20の次元を有し、Lはバリエーションの周りのランキング配列の長さである(我々の事例では、Lは51個のアミノ酸を表す)。

#### 【0476】

事前訓練された3状態二次構造および3状態溶媒接触性ネットワークへの入力のために、

やはり長さが51であり深さが20である、すべての99種の脊椎動物に対する複数配列アラインメントから生成される単一のPFM行列を使用した。Protein DataBankからの既知の結晶構造についてネットワークを事前訓練した後で、二次構造および溶媒モデルに対する最後の2つの層が除去され(グローバル最大プリーング層および出力層)、以前の層の出力の51×40の形状が、病原性予測ネットワークに対する入力として使用された。パラメータを精密に調整するために、ネットワークの構造層を通じた逆伝播を許容した。

#### 【0477】

##### [半教師あり学習]

半教師あり学習アルゴリズムは、訓練プロセスにおいてラベリングされたインスタンスとラベリングされていないインスタンスの両方を使用するので、訓練に利用可能な少量のラベリングされたデータしかない完全教師あり (completely supervised) 学習アルゴリズムよりも高い性能を達成する分類器を生み出すことができる。半教師あり学習の背後にある原理は、ラベリングされたインスタンスだけを使用する教師ありモデルの予測能力を強化するために、ラベリングされていないデータ内の固有の知識を活用できるということであり、それにより半教師あり学習の潜在的な利益がもたらされる。少量のラベリングされたデータから教師あり分類器により学習されるモデルパラメータは、ラベリングされていないデータによって、より現実的な分布(これは検定データの分布によく似ている)に向かって導かれ得る。

10

#### 【0478】

バイオインフォマティクスにおいて広まっている別の課題はデータ不均衡の問題である。データ不均衡現象は、予測されるべきクラスのうちの1つがデータにおいて過小評価されているときに生じ、それは、そのクラスに属するインスタンスが稀であり(注目に値する事例)、または取得が難しいからである。皮肉なことに、少数派のクラスが通常は最も学習に重要であり、それはそれらが特別な事例と関連付けられ得るからである。

20

#### 【0479】

不均衡なデータ分布に対処するためのアルゴリズム手法は、分類器のアンサンブルに基づく。ラベリングされたデータの量が限られていることは、当然より弱い分類器につながるが、弱い分類器のアンサンブルはどのような単一の構成分類器の性能も超える傾向がある。その上、アンサンブルは通常、複数のモデルを学習することと関連付けられる労力およびコストの妥当性を立証する要因によって、単一の分類器から取得される予測の正確さを改善する。直観的には、いくつかの分類器を集約することはより優れた過剰適合の制御につながり、それは、個々の分類器の高い変動性を平均化することは分類器の過剰適合も平均化するからである。

30

#### 【0480】

高い信頼性でラベリングされた病原性バリエーションの適切なサイズのデータセットがなかったため、半教師あり学習の戦略を追求した。ClinVarデータベースは300000個のエントリを有するが、有意性が不確かなバリエーションを除去した後で、病原性の解釈が矛盾しないミスセンスバリエーションは約42000個しか残らなかった。

#### 【0481】

体系的な検討により、これらのエントリがアノテートされた病原性を支持するのに十分な臨床的エビデンスを持っていないことも発見された。その上、人により精選されたデータベースの中のバリエーションの大半は、遺伝子の非常に小さい集合の中にある傾向があり、これにより、それらのバリエーションが、一般的なヒトバリエーションまたはチンパンジーとヒトとで固定された置換を使用してゲノムワイドで確認された良性訓練データセットの中のバリエーションと一致しなくなる。データセットがどれだけ異なるように確認されたかを考慮すると、人により精選されたバリエーションを病原性セットとして、およびゲノム全体の一般的なバリエーションを良性セットとして用いて、教師あり学習モデル訓練することは、重大なバイアスをもたらす可能性が高い。

40

#### 【0482】

バイアスを除去するように注意深く照合された、ラベリングされた良性バリエーションのセ

50

ットとバリエーションのラベリングされていないセットとを区別するように、深層学習ネットワークを訓練した。一実装形態によれば、385236個のラベリングされた良性バリエーションのセットは、ExAC/gnomADデータベースからの一般的なヒットバリエーション(>0.1%アレル頻度)およびヒット以外の霊長類の6つの種からのバリエーションを含んでいた。

#### 【0483】

(変異率、遺伝的浮動、および遺伝子変換を考慮するために)トリヌクレオチドコンテキストでの良性バリエーションとの一致を条件とし、バリエーション確認に対するアラインメント可能性およびシーケンスカバレッジの影響を調整して、ラベリングされていないバリエーションのセットをサンプリングした。ラベリングされていないバリエーションの数はラベリングされた良性バリエーションを大きく超えるので、ラベリングされた良性バリエーションの同じセットおよびラベリングされていないバリエーションの8つのランダムにサンプリングされたセットを使用する8つのモデルを訓練し、それらの予測の平均をとることによって、コンセンサス予測を得た。

10

#### 【0484】

半教師あり学習を選ぶ動機は、人により精選されたバリエーションデータベースが、信頼できずノイズが多く、特に、信頼性のある病原性バリエーションを欠いているということにある。gnomADからの一般的なヒットバリエーションおよび霊長類バリエーションから、信頼性のある良性バリエーションのセットを得た。病原性バリエーションについては、未知のバリエーションのセット(臨床的な有意性がアノテートされていないVUSバリエーション)からの病原性バリエーションを最初にサンプリングすることに対して、反復的バランスサンプリング(iterative balanced sampling)手法を採用した。

20

#### 【0485】

サンプリングバイアスを下げるために、良性訓練バリエーションの同じセットおよび病原性バリエーションの8つの異なるセットを使用する、8つのモデルのアンサンブルを訓練した。最初に、病原性バリエーションを表現するために未知のバリエーションをランダムにサンプリングした。次に、反復的に、モデルのアンサンブルが、前の訓練サイクルに関与していなかった未知のバリエーションのセットをスコアリングするために使用される。次いで、上位のスコアリングされた病原性バリエーションが、前のサイクルにおけるランダムな未知のバリエーションの5%を置き換えるために取得される。必要とされるよりも25%多くの上位にスコアリングされた病原性バリエーションを保持したので、8つのモデルに対するランダム性を高める未知のバリエーションを置き換えるために、スコアリングされた病原性バリエーションの8つの異なるセットをサンプリングできることに留意されたい。次いで、新しい病原性訓練セットが形成され、新しい訓練のサイクルが実行される。このプロセスは、最初のランダムにサンプリングされた未知のバリエーションが、アンサンブルモデルによって予測される信頼性の高い病原性バリエーションによりすべて置き換えられるまで繰り返される。図42は、反復的バランスサンプリングプロセスの例示である。

30

#### 【0486】

[ 良性訓練セットおよび未知の訓練セットのバランスをとること ]

良性バリエーションと一致している未知のバリエーションのサンプリング方式は、我々のモデル訓練のバイアスを低減するのに有用である。未知のバリエーションがランダムにサンプリングされるとき、深層学習モデルはしばしば、偏った情報を抽出して自明解を提示する。たとえば、アミノ酸置換K Mが良性バリエーションより未知のバリエーションにおいて頻繁に発生する場合、深層学習モデルはK Mの置換を常に病原性として分類する傾向がある。したがって、2つの訓練セットの間でアミノ酸配列の分布のバランスをとることが重要である。

40

#### 【0487】

CpGトランジションのようなより変異可能性の高いクラスは、一般的な良性バリエーションにおいて大きな表現バイアスを有する。他の霊長類からのオーソログバリエーションもヒットの変異率に従い、良性訓練セット全体における変異可能性の高いクラスのエンリッチメントを示唆する。未知のバリエーションのサンプリング手順があまり制御されておらず、バランスがとれていない場合、深層学習モデルは、トランスポージョンまたは非CpGトランジシ

50

ョンなどのより出現しないクラスと比較して、CpGトランジションを良性として分類する可能性がより高い傾向がある。

【0488】

深層学習モデルが自明で非生物学的な解に収束するのを防ぐために、良性バリエーションおよび未知のバリエーションのトリヌクレオチドコンテキストのバランスをとることを考える。トリヌクレオチドは、バリエーションの前の塩基、バリエーションの基準塩基、およびバリエーションの後の塩基によって形成される。そして、バリエーションの基準塩基は他の3つのヌクレオチドへと変更され得る。全体で、64×3個のトリヌクレオチドコンテキストがある。

【0489】

[反復的バランスサンプリング]

10

サイクル1

各トリヌクレオチドコンテキストに対する良性バリエーションの厳密な数と一致するように未知のバリエーションをサンプリングした。言い換えると、最初のサイクルにおいて、バリエーションのトリヌクレオチドコンテキストに関して良性訓練セットおよび病原性訓練セットを鏡写しにした。そのようなサンプリング方法の背後にある直観は、良性セットと未知のセットの間で変異率が同一であるバリエーションの等しい表現があるということである。このことは、モデルが変異率に基づいて自明解に収束するのを防ぐ。

【0490】

サイクル2～サイクル20

サイクル2に対して、サイクル1からの訓練されたモデルを適用してサイクル1に關与していない未知のバリエーションのセットをスコアリングし、上位の予測される病原性バリエーションで未知のバリエーションの5%を置き換えた。このセットは純粋にモデルによって生成され、このセットの中のトリヌクレオチドコンテキストに対するバランスは適用しなかった。訓練に必要な未知のバリエーションの残りの95%は、良性バリエーションの中の各ヌクレオチドコンテキストのカウントの95%となるようにサンプリングされる。

20

【0491】

直観的には、サイクル1は完全に一致した訓練セットを使用するので、上位の予測される病原性バリエーションはどのような変異率のバイアスも伴わずに生成される。したがって、このセットにおいてはどのようなバイアスも考慮する必要はない。データの残りの95%は依然として、モデルが自明解へ収束するのを防ぐために、トリヌクレオチドコンテキスト変異率のために制御される。

30

【0492】

各サイクルに対して、置き換えられた未知のバリエーションの百分率は5%ずつ上昇する。サイクル3では、サイクル3のモデルから上位の予測される病原性バリエーションで未知のバリエーションの5%を置き換えた。累積すると、病原性バリエーションの割合は10%に上昇し、トリヌクレオチドコンテキストと鏡写しにされた未知のバリエーションの割合は90%に下落する。サンプリングプロセスは残りのサイクルに対しても同様である。

【0493】

サイクル21

最後のサイクルであるサイクル21では、病原性訓練セット全体が、純粋に深層学習モデルから予測される上位の病原性バリエーションからなる。各サイクルにおいて変異率のバイアスを明確に考慮してきたので、病原性バリエーションは、訓練データとして使用するのに信頼性が高く、変異率のバイアスの影響を受けていない。したがって、訓練の最後のサイクルは、病原性予測のための最後の深層学習モデルを生み出す。

40

【0494】

[ラベリングされた良性訓練セットとラベリングされていない訓練セットの照合]

ラベリングされていないバリエーションのバランスサンプリングが、バリエーションの有害性に関連しないバイアスを除去するのに決定的に重要である。混乱をもたらす影響の適切な制御がないと、深層学習は容易に、不注意にもたらされたバイアスを選択してクラスを区別することがある。一般的なヒトバリエーションは、CpGアイランド上のバリエーションなどの、変

50



異可能性の高いクラスからのバリエーションについてエンリッチされる傾向がある。同様に、霊長類多型はヒトの変異率にも従い、良性訓練セット全体における変異可能性の高いバリエーションのエンリッチメントを示唆する。ラベリングされていないバリエーションのサンプリング手順がよく制御されておらずバランスがとれていない場合、深層学習ネットワークは、バリエーションを分類するために変異率のバイアスに頼る傾向があるので、トランスバージョンまたは非CpGトランジションなどのより出現しないクラスと比較して、CpGトランジションを良性として分類する可能性がより高い。我々は、96個のトリヌクレオチドコンテキスト(上で論じられた)の各々において、ラベリングされた良性バリエーションと厳密に同じ数のラベリングされていないバリエーションをサンプリングした。

#### 【0495】

ラベリングされた良性データセットの中の霊長類バリエーションに対してラベリングされていないバリエーションを照合するとき、我々は、複数配列アラインメントにおいて霊長類の種がアラインメントされなかったヒトゲノムの領域から、ラベリングされていないバリエーションが選択されるのを認めなかった。それは、その場所においてその霊長類の種の中のバリエーションをコールすることが可能ではなかったからである。

#### 【0496】

96個のトリヌクレオチドコンテキストの各々の中で、霊長類バリエーションに対するシーケンシングカバレッジを訂正した。シーケンシングされるヒトの数が多いので、ヒト集団における一般的なバリエーションは、シーケンシングカバレッジが低いエリアにおいても、それらが十分に確認されるほど十分に頻りに観察される。同じことは霊長類バリエーションに当てはまらず、それは、少数の個体しかシーケンシングされなかったからである。ExAC/gnomADエクソンにおけるシーケンシングカバレッジに基づいて、ゲノムを10個のビンへと分割した。各ビンに対して、ラベリングされた良性データセットvsラベリングされていないデータセットにおける霊長類バリエーションの割合を測定した。シーケンシングカバレッジのみに基づいて、線形回帰を使用して、霊長類バリエーションがラベリングされた良性データセットのメンバーである確率を計算した(図24)。ラベリングされた良性データセットにおける霊長類バリエーションと照合すべきラベリングされていないバリエーションを選択するとき、回帰係数を使用してその場所におけるシーケンシングカバレッジに基づいてバリエーションをサンプリングする確率を重み付けた。

#### 【0497】

[ 良性バリエーションおよび未知のバリエーションの生成 ]

ヒト集団における一般的なバリエーション

最近の研究は、ヒト集団における一般的なバリエーションが全般に良性であることを実証している。一実装形態によれば、gnomADは、正規のコーディング領域内でマイナーアレル頻度(MAF)が0.1%以上である90958個の非同義SNPを提供する。フィルタを通過したバリエーションが保持される。インデルが除外される。開始コドンまたは終止コドンにおいて発生するバリエーション、ならびにタンパク質切断バリエーションが除去される。亜集団を精査すると、各亜集団内のMAFが0.1%以上であるミスセンスバリエーションの総数は、一実装形態によれば245360個まで増える。これらのバリエーションは、良性バリエーションの訓練セットの一部を形成する。

#### 【0498】

大型類人猿における一般的な多型

コーディング領域は高度に保存的であることが知られているので、多型が大型類人猿の集団において高い頻度で分離しているかどうかを仮定するのは簡単であり、多型は健康に対する軽度の影響も有し得る。大型類人猿ゲノムプロジェクトおよび他の研究からの、ボノボ、チンパンジー、ゴリラ、およびオランウータンの多型データは、dbSNPからのアカゲザルおよびマーモセットのSNPと統合された。

#### 【0499】

未知のバリエーションの生成

すべての潜在的なバリエーションが、正規のコーディング領域の各塩基場所から、その場所

10

20

30

40

50

におけるヌクレオチドを他の3つのヌクレオチドに置換することによって生成される。新しいコドンが形成され、その場所におけるアミノ酸の潜在的な変更につながる。同義変化はフィルタリングされる。

#### 【0500】

gnomADデータセットにおいて観察されるバリエーションが除去される。開始コドンまたは終止コドンにおいて発生するバリエーション、ならびに終止コドンを形成するバリエーションが除去される。複数の遺伝子アノテーションを持つSNPに対して、正規の遺伝子アノテーションが、SNPのアノテーションを表すために選択される。全体で、一実装形態によれば、68258623個の未知のバリエーションが生成される。

#### 【0501】

##### バリエーションの追加のフィルタリング

ヒトゲノムの一部の領域では、リードをアラインメントするのが難しいことが知られている。それらの領域を含めると、訓練データセットおよび検定データセットに混乱をもたらす影響を引き起こす。たとえば、高い選択圧を受ける領域は、多型の数が限られる傾向がある。一方、シーケンシングが難しい領域もより少数の多型を有する。我々のモデルへのそのような混乱をもたらす入力为了避免するために、gnomADによってシーケンシングされなかった遺伝子からのバリエーションを除去した。

#### 【0502】

一般に、良性バリエーションは、複数の種にわたって保存される傾向があるよくシーケンシングされた領域において発見される。一方で、未知のバリエーションは、いくつかのカバー率の低い領域を含むゲノム全体でランダムにサンプリングされる。これは、良性セットと未知のセットとの間に確認バイアスを引き起こす。バイアスを減らすために、gnomADにおいてリード深さが10未満であるバリエーションを除去した。すべての哺乳類の種にわたるランキング配列アラインメントにおいて10%を超える欠けているデータがあるすべてのバリエーションも除去した。

#### 【0503】

##### 妥当性確認および検定のためのデータ

病原性モデルの妥当性確認および検定のために、一実装形態によれば、妥当性確認および検定のために、それぞれ10000個の良性バリエーションの2つのセットを、良性バリエーションの大きいプールからランダムにサンプリングした。良性バリエーションの残りは、深層学習モデルを訓練するために使用される。これらのバリエーションは特に、方法間の公平な比較を確実にするためにオーソロガスな霊長類バリエーションからサンプリングされ、それは、一部の方法が一般的なヒトバリエーションについて訓練されるからである。一実装形態によれば、妥当性確認および検定のために別々に、10000個の未知のバリエーションの2つのセットをランダムにサンプリングした。192個のトリヌクレオチドコンテキストの各々の中の未知のバリエーションの数が、妥当性確認セットおよび検定セットに対するそれぞれの良性バリエーションの数と一致することを確実にする。

#### 【0504】

自閉症または発育不全障害(DDD)を持つ影響を受けている子供と、影響を受けていない兄弟のde novoバリエーションを使用して、臨床上の環境において複数の方法の性能を評価した。全体で、一実装形態によれば、DDDの症例群からの3821個のミスセンスde novoバリエーションがあり、自閉症の症例群からの2736個のミスセンスde novoバリエーションがある。一実装形態によれば、影響を受けていない兄弟について1231個のミスセンスde novoバリエーションがある。

#### 【0505】

##### [ 深層学習ネットワークアーキテクチャ ]

病原性予測ネットワークは、二次構造および溶媒接触性ネットワークを介して、5つの直接入力および2つの間接入力を受け取る。5つの直接入力、長さ51個のアミノ酸配列 × 深さ20(20個の異なるアミノ酸を符号化する)であり、バリエーションを伴わない基準ヒトアミノ酸配列(1a)と、バリエーションで置換された代替ヒトアミノ酸配列(1b)と、霊長類の種の複

10

20

30

40

50

数配列アラインメントからのPFM(1c)と、哺乳類の種の複数配列アラインメントからのPFM(1d)と、より遠縁の脊椎動物の種の複数配列アラインメントからのPFM(1e)とを備える。二次構造および溶媒接触性ネットワークは各々、複数配列アラインメント(1f)および(1g)からのPFMを入力として受け取り、主な病原性予測ネットワークへの入力として出力を提供する。二次構造および溶媒接触性ネットワークは、Protein DataBankのための既知のタンパク質結晶構造について事前訓練され、病原性モデル訓練の間の逆伝播を可能にする。

#### 【0506】

5つの直接入力チャンネルは、線形活性化を伴う40個のカーネルのアップサンプリング畳み込み層を通される。ヒト基準アミノ酸配列(1a)は、霊長類、哺乳類、および脊椎動物の複数配列アラインメントからのPFMと統合される(マージ1a)。同様に、ヒト代替アミノ酸配列(1b)は、霊長類、哺乳類、および脊椎動物の複数配列アラインメントからのPFMと統合される(マージ1b)。これは2つの並列な経路を作り出し、1つは基準配列のためのものであり、1つはバリエーションで置換された代替配列のためのものである。

10

#### 【0507】

基準チャンネルと代替チャンネルの両方の統合された特徴マップ(マージ1aおよびマージ1b)は、一連の6つの残差ブロックを通される(層2a~7a、マージ2aおよび層2b~7b、マージ2b)。残差ブロックの出力(マージ2aおよびマージ2b)は、サイズ(51,80)の特徴マップを形成するために一緒に連結され(マージ3a、マージ3b)、これは基準チャンネルと代替チャンネルからのデータを完全に混合する。次に、データは、セクション2.1において定義されるような2つの畳み込み層を各々含む一連の6つの残差ブロックを通じて(マージ3~9、層21、層34を除く層9~46)、または、1D畳み込みを通った後の1つおきの残差ブロックの出力を接続するスキップ接続を介して(層21、層37、層47)のいずれかで、ネットワークを並列に通るための2つの経路を有する。最終的に、統合された活性化(マージ10)が別の残差ブロックに供給される(層48~53、マージ11)。マージ11からの活性化は、フィルタサイズ1の1D畳み込みおよびシングモイド活性化に与えられ(層54)、次いで、ネットワークによるバリエーション病原性の予測を表す単一の値を選ぶグローバル最大プーリング層に通される。モデルの概略的な図示が図3および補足テーブル16に示されている(図4A、図4B、および図4C)。

20

#### 【0508】

##### [モデル概要]

バリエーションの病原性を予測するために、半教師あり深層畳み込みニューラルネットワーク(CNN)モデルを開発した。モデルへの入力特徴量は、フランキングバリエーションのタンパク質配列および保存プロファイルと、特定の遺伝子領域におけるミスセンスバリエーションの枯渇率とを含む。深層学習モデルによって二次構造および溶媒接触性へバリエーションによって引き起こされる変化を予測し、それを我々の病原性予測モデルへと統合した。モデルを訓練するために、ヒト垂集団の一般的なバリエーションからの良性バリエーションと、霊長類からのオーソログバリエーションとを生成した。しかしながら、病原性バリエーションに対する信頼性のある源が依然として欠けている。最初に、良性バリエーションおよび未知のバリエーションを用いてモデルを訓練し、次いで、半教師あり反復的バランスサンプリング(IBS)アルゴリズムを使用して、高い信頼性で予測される病原性バリエーションのセットで未知のバリエーションを徐々に置き換えた。最終的に、ヒトにおいて発育不全障害を引き起こすde novoバリエーションを良性のバリエーションから区別する際に、我々のモデルが既存の方法を上回ることを実証した。

30

40

#### 【0509】

##### [残差ブロックの採用]

図17は残差ブロックを示す。病原性予測の我々の深層学習モデルと、二次構造および溶媒接触性を予測するための深層学習モデルの両方が、において最初に示された残差ブロックの定義を採用する。残差ブロックの構造は以下の図において示される。入力層は、まずバッチ正規化され、非線形活性化「ReLU」がそれに続く。活性化は次いで1D畳み込み層に通される。1D畳み込み層からのこの中間出力は、再びバッチ正規化およびReLU活性化され、別の1D畳み込み層が後に続く。第2の1D畳み込みの終わりにおいて、その出力を元の出

50

力と統合する。そのようなアーキテクチャでは、入力はその元の状態に保たれ、残差接続はモデルの非線形活性化がない状態に保たれる。

#### 【0510】

膨張/拡張畳み込みは、少数の訓練可能なパラメータで大きな受容野を可能にする。膨張/拡張畳み込みは、膨張畳み込み率または拡張係数とも呼ばれる、何らかのステップを用いて入力値をスキップすることによって、カーネルがその長さより大きいエリアにわたって適用されるような畳み込みである。膨張/拡張畳み込みは、畳み込み演算が実行されるときにより長い間隔の隣接する入力エントリ(たとえば、ヌクレオチド、アミノ酸)が考慮されるように、畳み込みフィルタ/カーネルの要素間に間隔を追加する。これにより、入力に長距離のコンテキスト依存性を組み込むことが可能になる。膨張畳み込みは、隣り合うヌクレオチドが処理される際に再使用するために部分的な畳み込み計算結果を保存する。

10

#### 【0511】

##### [我々のモデルの新規性]

我々の方法は、3つの点でバリエーションの病原性を予測するための既存の方法と異なる。第1に、我々の方法は、半教師あり深層畳み込みニューラルネットワークの新規のアーキテクチャを採用する。第2に、信頼性のある良性バリエーションがgnomADからの一般的なヒトバリエーションおよび霊長類バリエーションから取得され、一方で、確実性の高い病原性訓練セットは、人により精選された同一のバリエーションデータベースを使用したモデルの循環的な訓練および検定を避けるために、反復的バランスサンプリングおよび訓練を通じて生成される。第3に、二次構造および溶媒接触性のための深層学習モデルは、我々の病原性モデルのアーキテクチャへと統合される。構造および溶媒モデルから得られる情報は、特定のアミノ酸残基に対するラベル予測に限定されない。むしろ、リードアウト層が構造および溶媒モデルから除去され、事前訓練されたモデルが病原性モデルと統合される。病原性モデルを訓練する間、事前訓練された構造および溶媒層はまた、誤差を最小限にするために逆伝播する。これは、事前訓練された構造および溶媒モデルが、病原性予測問題に集中することを助ける。

20

#### 【0512】

##### [二次構造および溶媒接触性モデルの訓練]

##### データ準備

タンパク質の3状態の二次構造および3状態の溶媒接触性を予測するために、深層畳み込みニューラルネットワークを訓練した。PDBからのタンパク質アノテーションが、モデルを訓練するために使用される。一実装形態によれば、配列プロファイルと25%を超える相同性を有する配列が除去される。一実装形態によれば、全体で、6293個のタンパク質配列が訓練のために使用され、392個が妥当性確認のために使用され、499個が検定のために使用される。

30

#### 【0513】

タンパク質に対する位置特異的スコアリング行列(PSSM)保存プロファイルが、UniRef90を探すためにE値の閾値0.001および3回の反復という条件でPSI-BLASTを実行することによって生成される。あらゆる未知のアミノ酸ならびにその二次構造は、空白として設定される。また、すべてのヒト遺伝子に対する同様のパラメータ設定を用いてPSI-BLASTを実行し、それらのPSSM保存プロファイルを収集する。これらの行列は、構造モデルを病原性予測へと統合するために使用される。タンパク質配列のアミノ酸は次いで、ワンホット符号化ベクトルへと変換される。そして、タンパク質配列およびPSSM行列はL×20の行列へと形状変更され、Lはタンパク質の長さである。二次構造に対する3つの予測されるラベルは、ヘリックス(H)、シート(B)、およびコイル(C)を含む。溶媒接触性に対する3つのラベルは、埋もれている(B)、中間(I)、および露出している(E)を含む。1つのラベルは1つのアミノ酸残基に対応する。ラベルは次元=3のワンホット符号化ベクトルとしてコーディングされる。

40

#### 【0514】

50

### モデルアーキテクチャおよび訓練

タンパク質の3状態の二次構造および3状態の溶媒接触性をそれぞれ予測するために、2つのエンドツーエンドの深層畳み込みニューラルネットワークモデルを訓練した。2つのモデルは同様の構成を有し、一方はタンパク質配列に対する、他方はタンパク質保存プロファイルに対する、2つの入力チャンネルを含む。各入力チャンネルは次元 $L \times 20$ を有し、 $L$ はタンパク質の長さを示す。

#### 【0515】

入力チャンネルの各々は、40個のカーネルおよび線形活性化を伴う1D畳み込み層(層1aおよび層1b)を通される。この層は入力次元を20から40にアップサンプリングするために使用される。モデル全体ですべての他の層が40個のカーネルを使用することに留意されたい。2つの層(1aおよび1b)の活性化は、40個の次元の各々にわたって値を加算することによって一緒に統合される(すなわち、マージモード=「加算」)。マージノードの出力は、1D畳み込みの単一の層(層2)に、続いて線形活性化に通される。

10

#### 【0516】

層2からの活性化は、上で定義されたような一連の9個の残差ブロック(層3~層11)を通される。層3の活性化は層4に供給され、層4の活性化は層5に供給され、以下同様である。3つごとの残差ブロック(層5、層8、および層11)の出力を直接合計するスキップ接続もある。統合された活性化は次いで、ReLU活性化とともに2つの1D畳み込み(層12および層13)へと供給される。層13からの活性化はソフトマックスリードアウト層に与えられる。ソフトマックスは、所与の入力に対する3クラスの出力の確率を計算する。

20

#### 【0517】

最良の二次構造モデルでは、1D畳み込みは1という膨張率を有する。溶媒接触性モデルに対して、最後の3つの残差ブロック(層9、層10、および層11)は、カーネルのカバレッジを上げるために2という膨張率を有する。タンパク質の二次構造は、近くにあるアミノ酸の相互作用に強く依存する。したがって、カーネルカバレッジがより高いモデルは、性能をわずかに改善する。一方で、溶媒接触性は、アミノ酸間の長距離の相互作用により影響を受ける。したがって、膨張畳み込みを使用する、カーネルのカバレッジが高いモデルに対して、その正確さはカバレッジが低いモデルの正確さより2%高い。

#### 【0518】

以下の表は、一実装形態による、3状態二次構造予測モデルの各層に対する活性化およびパラメータについての詳細を提供する。

30

#### 【0519】

【表 2】

層	タイプ	カーネルの数、 ウィンドウ サイズ	形状	膨張率	活性化
入力配列(層 1a)	畳み込み 1D	40,1	(L,40)	1	線形
入力 PSSM(層 1b)	畳み込み 1D	40,1	(L,40)	1	線形
マージ配列 +PSSM	マージ(モード= 連結)	-	(L,80)	-	-
層 2	畳み込み 1D	40,5	(L,40)	1	線形
層 3	畳み込み 1D	40,5	(L,40)	1	ReLU
層 4	畳み込み 1D	40,5	(L,40)	1	ReLU
層 5	畳み込み 1D	40,5	(L,40)	1	ReLU
層 6	畳み込み 1D	40,5	(L,40)	1	ReLU
層 7	畳み込み 1D	40,5	(L,40)	1	ReLU
層 8	畳み込み 1D	40,5	(L,40)	1	ReLU
層 9	畳み込み 1D	40,5	(L,40)	1	ReLU
層 10	畳み込み 1D	40,5	(L,40)	1	ReLU
層 11	畳み込み 1D	40,5	(L,40)	1	ReLU
マージ活性化	マージ-層 5、層 8 および層 11 モ ード=「加算」	-	(L,40)	-	-
層 12	畳み込み 1D	40,5	(L,40)	1	ReLU
層 13	畳み込み 1D	40,5	(L,40)	1	ReLU
出力層	畳み込み 1D	1,3	(L,3)	-	ソフトマックス

10

20

30

40

## 【 0 5 2 0 】

一実装形態による、溶媒接触性モデルの詳細が以下の表に示されている。

## 【 0 5 2 1 】

【表 3】

層	タイプ	カーネルの数、 ウィンドウ サイズ	形状	膨張率	活性化
入力配列(層 1a)	畳み込み 1D	40,1	(L,40)	1	線形
入力 PSSM(層 1b)	畳み込み 1D	40,1	(L,40)	1	線形
マージ配列 +PSSM	マージ(モード=連 結)	-	(L,80)	-	-
層 2	畳み込み 1D	40,5	(L,40)	1	線形
層 3	畳み込み 1D	40,5	(L,40)	1	ReLU
層 4	畳み込み 1D	40,5	(L,40)	1	ReLU
層 5	畳み込み 1D	40,5	(L,40)	1	ReLU
層 6	畳み込み 1D	40,5	(L,40)	1	ReLU
層 7	畳み込み 1D	40,5	(L,40)	1	ReLU
層 8	畳み込み 1D	40,5	(L,40)	1	ReLU
層 9	畳み込み 1D	40,5	(L,40)	2	ReLU
層 10	畳み込み 1D	40,5	(L,40)	2	ReLU
層 11	畳み込み 1D	40,5	(L,40)	2	ReLU
マージ活性化	マージ-層 5、層 8 および層 11 モ ード=「加算」	-	(L,40)	-	-
層 12	畳み込み 1D	40,5	(L,40)	1	ReLU
層 13	畳み込み 1D	40,5	(L,40)	1	ReLU
出力層	畳み込み 1D	1,3	(L,3)	-	ソフトマックス

10

20

30

40

## 【0522】

特定のアミノ酸残基の二次構造クラスは、最大の予測されるソフトマックス確率によって決定される。モデルは、逆伝播を最適化するためのADAM最適化器を使用してタンパク質配列全体に対して累積カテゴリクロスエントロピー損失関数を用いて訓練される。

## 【0523】

3状態二次構造予測モデルに対する最高の検定の正確さは80.32%であり、同様の訓練データセット上でDeepCNFモデルによって予測される最新の正確さと同様である。

## 【0524】

50

3状態溶媒接触性予測モデルに対する最高の検定の正確さは64.83%であり、同様の訓練データセット上でRaptorXによって予測される現在の最高の正確さと同様である。

【0525】

以下で説明されるように、事前訓練された3状態二次構造および溶媒接触性予測モデルを、我々の病原性予測モデルへと統合した。

【0526】

[バリエーションの病原性を予測するようにモデルを訓練すること]

病原性予測モデルの入力特徴量

上で論じられたように、病原性予測問題に対して、病原性モデルを訓練するための良性バリエーション訓練セットおよび未知のバリエーション訓練セットがある。各バリエーションに対して、我々のモデルに供給するために以下の入力特徴量を準備した。

【0527】

各バリエーションの第1の入力特徴量は、バリエーションの配列コンテキストの深層学習モデルを提供するための、バリエーションのランキングアミノ酸配列、すなわち、hg19の基準配列から得られたバリエーションの各側の25個のアミノ酸である。全体で、このランキング基準配列は、51個のアミノ酸の長さを有する。

【0528】

第2の特徴量は、バリエーションを引き起こした代替アミノ酸である。基準アミノ酸と代替アミノ酸のペアを直接提供する代わりに、代替ランキング配列をモデルに提供する。代替ランキング配列は、配列の中間の場所が基準アミノ酸ではなく代替アミノ酸を含むことを除き、第1の特徴量における基準ランキング配列と同じである。

【0529】

次いで、両方の配列が長さ51×20のワンホット符号化されたベクトルへと変換され、各アミノ酸は20個の0または1のベクトルによって表される。

【0530】

次いで、12種の霊長類のための1つの位置特重的重み行列(PWM)、霊長類を除く47種の哺乳類のための1つのPWM、および霊長類と哺乳類を除く40種の脊椎動物のための1つのPWMを含む、3つのPWMが、バリエーションに対する99種の脊椎動物の複数配列アラインメント(MSA)から生成される。各PWMはL×20という次元を有し、Lはバリエーションの周りのランキング配列の長さである(我々の場合ではLは51個のアミノ酸を表す)。これは、種の各カテゴリにおいて見られるアミノ酸のカウントを備える。

【0531】

psi blastから、51個のアミノ酸のバリエーションランキング配列に対するPSSM行列も生成する。これは、病原性予測のための、3状態二次構造予測モデルおよび溶媒接触性予測モデルの統合に使用される。

【0532】

基準配列(入力1)、代替配列(入力2)、霊長類のためのPWM行列(入力3)、哺乳類(入力4)、脊椎動物(入力5)、ならびに3状態二次構造および溶媒接触性モデルからの情報を用いて、病原性モデルを訓練する。

【0533】

[深層学習モデルの訓練]

図19は、深層学習モデルワークフローの概要を提供するブロック図である。病原性訓練モデルは、5つの直接入力および4つの間接入力を備える。5つの直接入力特徴量は、基準配列(1a)、代替配列(1b)、霊長類保存率(1c)、哺乳類保存率(1d)、および脊椎動物保存率(1e)を含む。間接入力は、基準配列ベース二次構造(1f)、代替配列ベース二次構造(1g)、基準配列ベース溶媒接触性(1h)、および代替配列ベース溶媒接触性(1i)を含む。

【0534】

間接入力1fおよび1gに対して、ソフトマックス層を除く、二次構造予測モデルの事前訓練された層をロードする。入力1fに対して、事前訓練された層は、バリエーションに対してPSI-BLASTによって生成されるPSSMとともにバリエーションに対するヒト基準配列に基づく。同

10

20

30

40

50



様に、入力1gに対して、二次構造予測モデルの事前訓練された層は、PSSM行列とともに入力としてのヒト代替配列に基づく。入力1hおよび1iは、それぞれ、バリエーションの基準配列および代替配列に対する溶媒接触性情報を含む同様の事前訓練されたチャンネルに対応する。

#### 【 0 5 3 5 】

5つの直接入力チャンネルは、線形活性化を伴う40個のカーネルのアップサンプリング畳み込み層を通される。層1a、層1c、層1d、および層1eは、40個の特徴量の次元にわたって加算された値と統合されて、層2aを作り出す。言い換えると、基準配列の特徴マップは、3つのタイプの保存率特徴マップと統合される。同様に、1b、1c、1d、および1eは、40個の特徴量の次元にわたって加算された値と統合されて、層2bを生成する。すなわち、代替配列の特徴量は、3つのタイプの保存率特徴量と統合される。

10

#### 【 0 5 3 6 】

層2aおよび2bは、ReLUの活性化を用いてバッチ正規化され、各々フィルタサイズ40の1D畳み込み層に通される(3aおよび3b)。層3aおよび層3bの出力は、互いに連結された特徴マップを伴う1f、1g、1h、および1iと統合される。言い換えると、保存プロファイルを伴う基準配列の特徴マップ、および保存プロファイルを伴う代替配列は、基準配列および代替配列の二次構造特徴マップ、ならびに基準配列および代替配列の溶媒接触性特徴マップと統合される(層4)。

#### 【 0 5 3 7 】

層4の出力は、6つの残差ブロック(層5、層6、層7、層8、層9、層10)を通される。最後の3つの残差ブロックは、カーネルにより高いカバレッジを与えるために、1D畳み込みに対して2という膨張率を有する。層10の出力は、フィルタサイズ1の1D畳み込みおよび活性化シグモイドを通される(層11)。層11の出力は、バリエーションに対する単一の値を選ぶグローバル最大プールを通される。この値はバリエーションの病原性を表す。病原性予測モデルの一実装形態の詳細が以下の表に示される。

20

#### 【 0 5 3 8 】

【表 4 A】

層	タイプ	カーネルの数 、ウィンドウ サイズ	形状	膨張率	活性化
基準配列(1a)	畳み込み 1D	40,1	(51,40)	1	線形
代替配列(1b)	畳み込み 1D	40,1	(51,40)	1	線形
霊長類保存率 (1c)	畳み込み 1D	40,1	(51,40)	1	線形
哺乳類保存率 (1d)	畳み込み 1D	40,1	(51,40)	1	線形
脊椎動物保存 率(1e)	畳み込み 1D	40,1	(51, 40)	1	線形
基準配列ベー ス二次構造(1f)	入力層	-	(51,40)	-	-
代替配列ベー ス二次構造 (1g)	入力層	-	(51,40)	-	-
基準配列ベー ス溶媒接触性 (1h)	入力層	-	(51,40)	-	-
代替配列ベー ス溶媒接触性 (1i)	入力層	-	(51,40)	-	-
基準配列マー ジ(2a)	マージ(モード=加 算)(1a,1c,1d,1e)	-	(51,40)	-	-

10

20

30

40

【表 4 B】

(【表 4 A】の続き)

代替配列マージ(2b)	マージ(モード=加算)(1b,1c,1d,1e)	-	(51,40)	-	-
3a	畳み込み 1D(2a)	40,5	(51,40)	1	ReLU
3b	畳み込み 1D(2b)	40,5	(51,40)	1	ReLU
4	マージ(モード=連結) (3a,3b,1f,1g,1h,1i)	-	(51,240)	-	-
5	畳み込み 1D	40,5	(51,40)	1	ReLU
6	畳み込み 1D	40,5	(51,40)	1	ReLU
7	畳み込み 1D	40,5	(51,40)	1	ReLU
8	畳み込み 1D	40,5	(51,40)	1	ReLU
9	畳み込み 1D	40,5	(51,40)	2	ReLU
10	畳み込み 1D	40,5	(51,40)	2	ReLU
11	畳み込み 1D	40,1	(51,1)	2	シグモイド
12	グローバル最大プーリング	-	1	-	-
出力層	活性化層	-	1	-	シグモイド

10

20

30

40

【 0 5 4 0】

[ アンサンプル ]

一実装形態では、我々の方法の各サイクルに対して、同じ良性データセットおよび8つの異なる未知のデータセットで訓練する8つの異なるモデルを実行し、8つのモデルにわたって評価データセットの予測を平均した。未知のパリアントの複数のランダムにサンプリングされたセットがモデルに提示されると、サンプリングバイアスを減らしてよく制御することができる。

【 0 5 4 1】

さらに、アンサンプルのアンサンプルという手法の採用が、我々の評価データセットに対する我々のモデルの性能を改善することができる。CADDは、10個のモデルのアンサンプル

50

ルを使用して、バリエーションをスコアリングするために10個すべてのモデルにわたる平均スコアを得る。ここで、同様のアンサンブル手法を使用することを試みた。1つのアンサンブルを使用して結果のベンチマークをとり、次いで、アンサンブルの数を増やして性能の向上を評価した。各アンサンブルは、同じ良性データセットおよび8つの異なる未知のデータセットで訓練する8つのモデルを有することに留意されたい。異なるアンサンブルに対して、乱数生成器のシード値は別個であるので、ランダムなバリエーションのセットが互いに異なるように導かれる。

【0542】

—実装形態による詳細な結果が以下の表に示される。

【0543】

10

【表5】

アンサンブルの数	DDD データセット についての- $\log(pvalue)$ (アンサン ブルの平均の平均)	DDD データセットについての- $\log(pvalue)$ (アンサンブルの平均の 中央値)
1	3.4e-27	3.4e-27
5	2.74e-29	3.8e-29
10	2.98e-29	2.55e-29
15	4.06e-29	3.88e-29
20	3.116e-29	3.05e-29
25	3.77e-29	3.31e-29
30	3.81e-29	3.34e-29

20

30

【0544】

1つのアンサンブルと比較して、5つのアンサンブルおよび10個のアンサンブルが、DDD データセットを使用して評価されるときにより大きなp値を生み出した。しかし、アンサンブルの数を増やすことでさらに性能が向上することはなく、アンサンブルに対する飽和を示している。アンサンブルは、多様な未知のバリエーションを用いてサンプリングバイアスを減らす。しかしながら、良性クラスと病原性クラスとで192個のトリヌクレオチドコンテキストを照合することも必要とされ、これは我々のサンプリング空間をかなり制限し、早い飽和につながる。アンサンブルのアンサンブルという手法は、モデル性能を大きく改善し、モデルについての我々の理解をさらに深めると結論付ける。

40

【0545】

[病原性モデルを訓練するための早期打ち切り]

信頼性のあるアノテートされた病原性バリエーションサンプルが欠けているので、モデル訓練のための打ち切り基準を定義するのは困難である。モデル評価における病原性バリエ

50

トの使用を避けるために、一実装形態では、オーソロガスな霊長類からの10000個の良性妥当性確認バリエーションと、10000個のトリヌクレオチドコンテキストが照合された未知のバリエーションとを使用した。モデルの各エポックを訓練した後、良性妥当性確認バリエーションおよび未知の妥当性確認バリエーションを評価した。妥当性確認バリエーションセットの両方の確率分布の差を評価するために、ウィルコクソン順位和検定を使用した。

#### 【0546】

検定のp値は、モデルが良性バリエーションを未知のバリエーションのセットから区別する能力の向上に伴って、より大きくなる。モデル訓練の任意の5つの連続するエポックの間に、2つの分布を区別するモデルの能力に改善が観察されない場合、訓練を打ち切る。

#### 【0547】

事前に、10000個の保留された霊長類バリエーションの2つの別個のセットを訓練から除外し、我々はこれらのセットを妥当性確認セットおよび検定セットと名付けた。モデル訓練の間の早期打ち切りを評価するために、10000個の保留された霊長類バリエーションおよび10000個のトリヌクレオチドコンテキストについて照合されたラベリングされていないバリエーションの妥当性確認セットを使用した。各訓練エポックの後で、ラベリングされた良性妥当性確認セットおよびラベリングされていない一致した対照群におけるバリエーションを区別するための、深層ニューラルネットワークの能力を評価し、ウィルコクソン順位和検定を使用して予測されるスコアの分布の差を測った。5つの連続的な訓練エポックの後でさらなる改善が観察されないと、過剰適合を防ぐために訓練を打ち切った。

#### 【0548】

[分類器の性能のベンチマーキング]

1つは一般的なヒトバリエーションのみを用いて訓練され、1つは一般的なヒトバリエーションと霊長類バリエーションの両方を含む良性とラベリングされた完全なデータセットを用いて訓練された、2つのバージョンの深層学習ネットワークの分類の正確さを、以下の分類器、すなわちSIFT、PolyPhen-2、CADD、REVEL、M-CAP、LRT、MutationTaster、MutationAssessor、FATHMM、PROVEAN、VEST3、MetaSVM、MetaLR、MutPred、DANN、FATHMM-MKL\_coding、Eigen、GenoCanyon、およびGERP++13,32-48に加えて評価した。他の分類器の各々のスコアを得るために、dbNSFP 49(<https://sites.google.com/site/jpopgen/dbNSFP>)からすべてのミスセンスバリエーションに対するスコアをダウンロードし、10000個の保留された霊長類バリエーション検定セット、およびDDD症例群vs対照群におけるde novoバリエーションについて方法のベンチマークをとった。本明細書に含めるものには、SIFT、PolyPhen-2、およびCADD、ならびにREVELを選択した。それは、SIFT、PolyPhen-2、およびCADDについては、それらが最も広く使用されている方法であるからであり、REVELについては、様々な評価モードにわたって、評価した20個の既存の分類器の中で最良のもの1つとして傑出していたからである。評価したすべての分類器の性能が図28Aにおいて提供される。

#### 【0549】

深層学習ネットワークの性能に対する、利用可能な訓練データサイズの影響を評価するために、385236個の霊長類バリエーションと一般的なヒトバリエーションのラベリングされた良性訓練セットからランダムにサンプリングすることによって、図6の各データ点において深層学習ネットワークを訓練した。分類器の性能におけるランダムなノイズを減らすために、初期パラメータ重みのランダムなインスタンス化を各回において使用してこの訓練手順を5回実行し、10000個の保留された霊長類バリエーションとDDD症例群vs対照群データセットの両方についての中央値の性能を図6に示した。偶然にも、385236個のラベリングされた良性バリエーションの完全なデータセットを用いた分類器の中央値の性能は、DDDデータセットについて本明細書の残りで使用したものよりわずかに高かった(ウィルコクソン順位和検定で $P < 10^{-28}$ ではなく $P < 10^{-29}$ )。各々の個別の霊長類の種からのバリエーションが分類の正確さに寄与しており、一方で各々の個別の哺乳類の種からのバリエーションが分類の正確さを下げることが示すために、一実装形態によれば、83546個のヒトバリエーションと、各種に対する一定数のランダムに選択されたバリエーションとを備える訓練データセットを使用して深層学習ネットワークを訓練した。一実装形態によれば、訓練セットに追加したバリエーション

10

20

30

40

50

の一定の数(23380)は、ミスセンスバリエーションの数が最小の種、すなわちボノボにおいて利用可能なバリエーションの総数である。ノイズを減らすために、訓練手順を再び5回繰り返し、分類器の中央値の性能を報告した。

【0550】

[モード評価]

—実装形態では、反復的バランスサンプリング手順に続いて、深層学習モデルを21回のサイクルにわたり訓練した。我々の分類器の性能を評価するために、2つのタイプの評価を実行した。2つの尺度で我々のモデルとPolyphen2、SIFT、およびCADDの比較も行い、臨床的なアノテーションに対する我々のモデルの適用の可能性を評価した。

【0551】

方法1: 良性検定セットの正確さ

—実装形態では、10000個の良性バリエーションおよび未知のバリエーションを、8つの異なる訓練されたモデルのアンサンプルを使用してそれらの予測される確率を計算することによって、評価した。上で言及された他の既存の方法によってスコアリングされる、それらの予測される確率も取得した。

【0552】

次いで、評価において使用される方法の各々に対する未知の検定バリエーションにわたる予測される確率の中央値を取得した。中央値スコアを使用して、方法の各々によって使用される良性バリエーションおよび病原性バリエーションのアノテーションに応じ、スコアが中央値の上または下である良性バリエーションの数を発見した。SIFT、CADD、および我々の方法は、病原性バリエーションを1として、良性バリエーションを0としてラベリングする。したがって、スコアが中央値の下である良性バリエーションの数をカウントした。Polyphenは反対のアノテーションを使用し、中央値を上回る良性バリエーションの数をカウントした。スコアが中央値の上/下である良性バリエーションの数を良性バリエーションの総数で割った比は、良性バリエーションの予測の正確さを表す。

良性の正確さ=中央値を上回る(下回る\*)良性バリエーションの総数÷良性バリエーションの総数

【0553】

この評価方法の背後にある我々の理論は、gnomADにおけるバリエーションの選択圧力の分析に依拠している。gnomADにおけるシングルトンに対して、ミスセンスバリエーションと同義バリエーションの比は約2.26:1である。一方、gnomADにおける一般的なバリエーション(MAF>0.1%)では、ミスセンス対同義比は約1.06:1である。これは、ランダムな未知のバリエーションのセットから、約50%が自然選択によって排除されることが予想され、残りの50%は軽度であり集団において一般的になる傾向があることを示している。

【0554】

10

20

30

【表 6】

方法	保留された良性セットの正確さ
Polyphen	0.74
SIFT	0.69
CADD	0.77
我々のモデル	0.85

10

## 【0555】

上の表に示されるように、我々の方法の性能は2番目に良い方法であるCADDより8%を超えて優れている。これは、良性バリエーションを分類する我々のモデルの能力の大きな向上を示している。そのような実証は我々のモデルの能力を証明するが、以下の方法2は、臨床的な解釈のための臨床的なデータセットに対する我々のモデルの有用性を示す。

20

## 【0556】

## 方法2: 臨床的なデータセットの評価

一実装形態では、発育不全障害(DDD)症例群-対照群データセットを含む、臨床的なデータセットに対してこれらの病原性予測方法を評価した。DDDデータセットは、影響を受けている子供からの3821個のde novoミスセンスバリエーションおよび影響を受けていない兄弟からの1231個のde novoミスセンスバリエーションを備える。我々の仮説は、影響を受けている子供からのde novoバリエーションが影響を受けていない兄弟からのde novoバリエーションより有害である傾向があるというものである。

30

## 【0557】

臨床的な検定データセットは病原性バリエーションを明確にラベリングしないので、それらの方法の性能を測るために、(影響を受けている群および影響を受けていない群からの)de novoバリエーションの2つのセットの分離を使用した。de novoバリエーションのこれらの2つのセットの分離がどれだけ優れているかを評価するために、ウィルコクソン順位和検定を適用した。

## 【0558】

【表7】

方法	DDD データセットについての log10(p 値)
Polyphen	15.02
SIFT	13.52
CADD	13.83
DL	28.35

10

## 【0559】

上の表によれば、我々の半教師あり深層学習モデルは、de novoバリエーションの影響を受けているセットと影響を受けていないセットとを区別する際の性能がはるかに高い。これは、我々のモデルが臨床的な解釈に対して既存の方法より適切であることを示している。これはまた、ゲノム配列および保存プロファイルから特徴量を抽出するという全般的な手法が、人により精選されたデータセットに基づく人が加工した特徴量より優れていることを確認する。

20

## 【0560】

[10000個の霊長類バリエーションの保留された検定セットについての良性予測の正確さ]

深層学習ネットワークならびに他の20個の分類器のベンチマークをとるために、検定データセットの中の10000個の保留された霊長類バリエーションを使用した。異なる分類器は大きく変動するスコア分布を有していたので、各分類器に対する50パーセントイル閾値を特定するために、トリヌクレオチドコンテキストにより検定セットと照合された、10000個のランダムに選択されたラベリングされていないバリエーションを使用した。方法間の公平な比較を確実にするために、その分類器に対して50パーセントイルの閾値で良性であると分類された、10000個の保留された霊長類バリエーション検定セットの中のバリエーションの割合について、各分類器のベンチマークをとった。

30

## 【0561】

良性バリエーションを特定するために50パーセントイルを使用することの背後にある我々の理論は、ExAC/gnomADデータセットの中のミスセンスバリエーションに対して観察される選択圧力に基づく。シングルTONアレル頻度で発生するバリエーションでは、ミスセンス:同義比は~2.2:1であるが、一般的なバリエーション(>0.1%アレル頻度)では、ミスセンス:同義比は約1.06:1である。これは、ミスセンスバリエーションの約50%が一般的なアレル頻度では自然選択により排除されることが予想され、残りの50%が遺伝子的浮動を介して集団において一般的になる可能性を有するのに十分軽度であることを示している。

40

## 【0562】

分類器の各々に対して、50パーセントイル閾値を使用して良性であるものとして予測される保留された霊長類検定バリエーションの割合が示されている(図28Aおよび補足テーブル17(図34))。

## 【0563】

[DDD研究からのde novoバリエーションの分析]

DDDの影響を受けている個人におけるde novoミスセンスバリエーションと、影響を受けてい

50



ない兄弟の対照群におけるde novoミスセンスバリエーションとを区別する能力について、分類方法のベンチマークをとった。各分類器に対して、2つの分布に対する予測スコア間の差のウィルコクソン順位和検定からのp値を報告した(図28Bおよび図28Cおよび補足テーブル17(図34))。

#### 【0564】

モデルの性能を分析するための我々の2つの尺度が異なる源および方法から導かれると仮定し、2つの異なる尺度についての分類器の性能が相関していたかどうかを検定した。実際に、我々はこれらの2つの尺度が相関していたことを発見し、保留された霊長類検定セットに対する良性分類の正確さと、DDD症例群vs対照群におけるde novoミスセンスバリエーションに対するウィルコクソン順位和検定のp値との間で、 $\text{spearman } \rho = 0.57 (P < 0.01)$ であった。これは、分類器のベンチマークをとることについて、保留された霊長類検定セットの正確さと、DDD症例群vs対照群のp値との間に、良い一致があることを示す(図30A)。

10

#### 【0565】

さらに、深層学習ネットワークが疾患と関連付けられる遺伝子の発見を助け得るかどうかを検定した。観察されたde novo変異の数をヌル変異モデルのもとで予想される数と比較することによって、遺伝子におけるde novo変異のエンリッチメントを検定した。

#### 【0566】

すべてのミスセンスde novo変異からの結果と、スコアが0.803より大きいミスセンス変異からの結果とを比較して、深層学習ネットワークの性能を調査した。すべてのミスセンスde novoを検定することにはデフォルトのミスセンス率を使用した。フィルタリングされたミスセンスde novoを検定することにはスコアが0.803より大きいサイトから計算されたミスセンス変異率を使用した。各遺伝子には4つの検定が必要であり、1つはタンパク質切断エンリッチメントを検定し、1つはタンパク質を変化させるde novo変異のエンリッチメントを検定し、両方がDDDコホートだけに対して、および神経発達トリオシーケンシングコホートのより大きなメタ分析に対して検定される。タンパク質を変化させるde novo変異のエンリッチメントは、コーディング配列内のミスセンスde novo変異のクラスタリングの検定とともに、Fisherの方法によって合成された(補足テーブル20および21)。各遺伝子に対するp値は4回の検定の最小値から取られ、ゲノムワイド有意性が $P < 6.757 \times 10^{-7}$ と決定された( $\alpha = 0.05$ 、18500個の遺伝子、4回の検定)。

20

#### 【0567】

[605個のDDD疾患関連遺伝子内での受信者操作特性曲線および分類の正確さの計算]

深層学習ネットワークが本当に同じ遺伝子内の病原性バリエーションと良性バリエーションとを区別していたかどうかを検定するために、de novo優性遺伝モードを伴う遺伝子における病原性を優先するのではなく、DDDコホートにおいてp値が0.05未満である(de novoタンパク質切断変異のみを使用して計算される)神経発達疾患と関連付けられた605の遺伝子のセットを特定した(補足テーブル18)。DDDデータセットおよび対照群データセットにおいて605個の遺伝子の中のバリエーションの確率分布を分類器が分離する能力について、すべての分類器に対するウィルコクソン順位和のp値を報告する(図28Cおよび補足テーブル19(図35))。

30

#### 【0568】

605個の遺伝子のこのセット内で、変異率だけから予想されるものの3倍のde novoミスセンスバリエーションのエンリッチメント比を観察した。これは、影響を受けているDDD患者におけるde novoミスセンスバリエーションが、約67%の病原性バリエーションと33%のバックグラウンドバリエーションを備え、一方で、健康な対照群におけるde novoミスセンスバリエーションが、不完全な浸透の事例を除いて大半はバックグラウンドバリエーションからなることを示している。

40

#### 【0569】

病原性バリエーションと良性バリエーションを完璧に区別する分類器に対する最大の可能なAUCを計算するために、605個の遺伝子内の影響を受けている個人におけるde novoミスセンスバリエーションの67%だけが病原性であり、残りがバックグラウンドであったことを考慮した

50

。受信者操作特性曲線を構築するために、病原性であるものとしてのde novo DDDバリアントの分類を真陽性のコールとして扱い、病原性であるものとしての健康な対照群におけるde novoバリアントの分類を偽陽性のコールとして扱った。したがって、完璧な分類器は、DDD患者におけるde novoバリアントの67%を真陽性として、DDD患者におけるde novoバリアントの33%を偽陽性として、対照群におけるde novoバリアントの100%を真陰性として分類する。受信者操作特性曲線の視覚化は、プロットの(0%,0%)の角および(100%,100%)の角へ直線によって接続された、真陽性率が67%であり偽陽性率が0%である単一の点を示すだけであり、良性変異と病原性変異の完璧な区別についての分類器に対する0.837という最大AUCが得られる(図30Bおよび補足テーブル19(図35))。

【0570】

合成されたDDDのデータセットおよび健康な対照群のデータセットにおける605個の遺伝子内での病原性バリアントの予想される割合を推定することによって、病原性バリアントと良性バリアントとをバイナリ閾値で分離するための深層学習ネットワークの分類の正確さを計算した。DDDデータセットは、予想を超える249個のde novoミスセンスバリアントの余剰を伴う379個のde novoバリアントを含み、対照群データセットは65個のde novoバリアントを含んでいたため、全体で444個のバリアントのうち249個の病原性バリアントを予測した(図22A)。この予想される比率に従って、444個のde novoミスセンスバリアントを良性または病原性カテゴリへと分離した各分類器に対する閾値を選択し、これをバイナリカットオフとして使用して各分類器の正確さを評価した。我々の深層学習モデルに対して、この閾値は、真陽性率が65%、偽陽性率が14%である0.803以上のカットオフにおいて達成された。DDDの個人における約33%のバックグラウンドバリアントの存在について調整された分類の正確さを計算するために、バックグラウンドであったde novo DDDバリアントの33%が、健康な対照群において観察したのと同じ偽陽性率で分類されるであろうと仮定した。これは、DDDデータセットにおける真陽性分類イベントの $14\% \times 0.33 = 4.6\%$ が実際にはバックグラウンドバリアントからの偽陽性であることに対応する。深層学習ネットワークに対する調整された真陽性率を、 $(65\% - 4.6\%) / 67\% = 90\%$ と推定する。深層学習ネットワークに対しては88%である、真陽性率および真陰性率の平均を報告する(図30Cおよび補足テーブル19(図35))。この推定は、神経発達障害において不完全な浸透が広く見られることにより、分類器の真の正確さを過小評価する可能性が高い。

【0571】

[ClinVar分類の正確さ]

既存の分類器の大半はClinVar上で訓練される。ClinVar上で直接訓練しない分類器も、ClinVar上で訓練される分類器からの予測スコアを使用することによって影響を受けることがある。加えて、一般的なヒトバリアントは良性のClinVarの結果に対して高度にエンリッチされ、それは、アレル頻度が、良性の結果をバリアントに割り当てるための基準の一部であるからである。

【0572】

我々は、ClinVarデータセットを分析に適したものにするために、2017年に追加されたClinVarバリアントだけを使用することによって、ClinVarデータセットにおける循環性を最小限にすることを試みた。それは、他の分類方法がその前の年に公開されたからである。2017年のClinVarバリアントの中でも、ExACにおいて一般的なアレル頻度(>0.1%)で存在するあらゆるバリアント、またはHGMD、LSDB、もしくはUniprotにおいて存在するあらゆるバリアントを除外した。すべてのそのようなバリアントをフィルタリングした後で、および有意性が不確かであり矛盾するアノテーションを伴うバリアントを除外した後で、ClinVarにおいて良性アノテーションを伴う177個のバリアントおよび病原性アノテーションを伴う969個のバリアントが残った。

【0573】

深層学習ネットワークと既存の方法の両方を使用してすべてのClinVarバリアントをスコアリングした。このデータセット内の良性バリアントおよび病原性バリアントの観察される比率に従って、ClinVarバリアントを良性カテゴリまたは病原性カテゴリに分離した

10

20

30

40

50

各分類器に対する閾値を選択し、これをバイナリカットオフとして使用して各分類器の正確さを評価した。各分類器に対する真陽性率および真陰性率の平均を報告する(図31Aおよび図31B)。ClinVarデータセットについての分類器の性能は、10000個の保留された霊長類バリエーションに対する分類の正確さについての分類器の性能、またはDDD症例群vs対照群データセットに対するウィルコクソン順位和のp値についての分類器の性能のいずれとも大きく相関しなかった(図31Aおよび図31B)。

#### 【0574】

我々は、既存の分類器は専門家の行動を正確にモデル化しているが、人の経験則は経験的なデータにおいて病原性変異と良性変異とを区別するのに完全に最適ではないことがあるという仮説を立てている。1つのそのような例はGranthamスコアであり、これはアミノ酸置換の相同性または非相同性を特徴付けるための距離の尺度を与える。完全なClinVarデータセット(約42000個のバリエーション)内の病原性バリエーションおよび良性バリエーションに対する平均Granthamスコアを計算し、これを605個の遺伝子内の影響を受けているDDDの個人および影響を受けていない個人におけるde novoバリエーションに対する平均Granthamスコアと比較した。影響を受けているDDDの個人における約33%のバックグラウンドバリエーションの存在を訂正するために、DDD症例群vs対照群の間のGranthamスコアの差を50%増大させたが、それでもこれは、ClinVarにおける病原性バリエーションと良性バリエーションとの差より小さかった。1つの可能性は、専門家が、アミノ酸置換距離などの、測定しやすい尺度を重視しすぎている一方で、専門家にとって定量化がより難しいタンパク質構造などの要因を軽視しているということである。

#### 【0575】

##### [ 深層学習モデルの解釈 ]

機械学習アルゴリズムが問題を解く手段を理解するのは難しいことが多い。バリエーションの病原性を予測するために深層学習ネットワークが学習して抽出した特徴量を理解するために、深層学習ネットワークの初期層を視覚化した。事前訓練された3状態二次構造予測モデルの最初の3つの層(2つのアップサンプリング層とそれに続く第1の畳み込み層)内での異なるアミノ酸に対する相関係数を計算し、BLOSUM62行列またはGrantham距離と非常に似た特徴量を畳み込み層の重みが学習することを示した。

#### 【0576】

異なるアミノ酸の間の相関係数を計算するために、二次構造モデルにおいて3つのアップサンプリング層(層1a、層1b、および層1c)が前にある第1の畳み込み層の重みから始めた。3つの層の間の行列乗算を実行し、次元が(20,5,40)である行列を得た。ここで、20はアミノ酸の数であり、5は畳み込み層のウィンドウサイズであり、40はカーネルの数である。最後の2つの次元を平坦化することによって次元(20,200)を有するように行列を形状変更し、20個のアミノ酸の各々に対して作用する重みが長さ200のベクトルとして表されるような行列を得た。20個のアミノ酸間の相関行列を計算した。各次元が各アミノ酸を表すので、相関係数行列を計算することによって、深層学習ネットワークが訓練データから学習したことに基づいて、アミノ酸間の相関と、深層学習ネットワークに対してアミノ酸がどれだけ類似しているように見えるかということとを計算している。相関係数行列の視覚化が図27に示されており(アミノ酸はBLOSUM62の行列順序でソートされている)、疎水性アミノ酸(メチオニン、イソロイシン、ロイシン、バリン、フェニルアラニン、チロシン、トリプトファン)および親水性アミノ酸(アスパラギン、アスパラギン酸、グルタミン酸、グルタミン、アルギニン、およびリジン)を備える2つの顕著なクラスターを示す。これらの初期層の出力はより後の層の入力になり、深層学習ネットワークがデータのますます複雑な階層的表現を構築することを可能にする。

#### 【0577】

ニューラルネットワークが予測において使用するアミノ酸配列のウィンドウを示すために、5000個のランダムに選択されたバリエーション内のおよびその周辺の各場所を摂動させて、バリエーションに対する予測されるPrimateAIスコアに対する影響を観察した(図25B)。バリエーションの周りの各々の近くのアミノ酸場所(-25~+25)における入力を系統的に0にして、

ニューラルネットワークにより予測されるバリエーションの病原性の変化を測定し、5000個のバリエーションにわたってその変化の平均絶対値をプロットした。バリエーションの近くのアミノ酸が最も大きい影響を受けており、概ね対称的な分布で、バリエーションからの距離が長くなるにつれて徐々に影響が小さくなる。重要なことに、このモデルは、バリエーションの場所におけるアミノ酸だけに基づくのではなく、タンパク質モチーフを認識するために必要とされるであろうより広いウィンドウからの情報を使用することによって、予測を行う。タンパク質サブドメインが比較的小型のサイズであることと一致して、51個を超えるアミノ酸へとウィンドウのサイズを延長することが、さらに正確さを改善しないことを経験的に観察した。

#### 【0578】

深層学習分類器のアラインメントに対する感度を評価するために、バリエーション分類の正確さに対するアラインメントの深さの影響を次のように調査した。アラインメントにおける種の数に基づいてデータを5つのピンへと分割し、各ピンにおけるネットワークの正確さを評価した(図57)。トリヌクレオチドコンテクストに対して照合された(図21Dのように、しかし各ピンに対して別々に実行される)ランダムに選択された変異から、保留された良性変異のセットを分離する際のネットワークの正確さは、上位の3つのピンにおいて最も高く、下位の2つのピンにおいて顕著に弱いことを発見した。99種の脊椎動物の多種アラインメントは、11種のヒト以外の霊長類、50種の哺乳類、および38種の脊椎動物を備え、下位の2つのピンは、他の非霊長類の哺乳類からのまばらなアラインメント情報を有するタンパク質を表す。深層学習ネットワークは、アラインメント情報が霊長類および哺乳類全体に及ぶときにロバストかつ正確であり、より遠縁の脊椎動物からの保存情報はより重要性が低い。

#### 【0579】

##### [ 正規のコーディング領域の定義 ]

正規のコーディング領域を定義するために、コーディングDNA配列(CDS)領域(knownCanonical.exonNuc.fa.gz)に対するヒトとの99種の脊椎動物ゲノムの複数アラインメントがUCSCゲノムブラウザからダウンロードされた。ヒトについては、エクソンの座標はBuild hg19のもとにある。エクソンは統合されて遺伝子を形成する。常染色体上の遺伝子およびchrXが保持される。相同ではない遺伝子は除去され、相同な遺伝子のリストはNCBI ftp://ftp.ncbi.nih.gov/pub/HomoloGene/current/homologene.dataからダウンロードされた。複数の遺伝子アノテーションを伴うSNPに対しては、SNPのアノテーションを表すために最長の転写産物が選択される。

#### 【0580】

##### [ ヒト、類人猿、および哺乳類の多型データ ]

世界中の8つの亜集団からの123136人の個人の全エクソンシーケンシングデータを収集した、最近の大規模な研究であるgenome Aggregation Database(gnomAD)から、ヒトエクソン多型データをダウンロードした。そして、フィルタを通過し正規のコーディング領域に該当するバリエーションを抽出した。

#### 【0581】

大型類人猿ゲノムシーケンシングプロジェクトは、24体のチンパンジー、13体のボノボ、27体のゴリラ、および10体のオランウータン(5体のスマトラオランウータンおよび5体のボルネオオランウータンを含む)の全ゲノムシーケンシングデータを提供する。チンパンジーおよびボノボについての研究は、追加の25体の類人猿のWGSを提供する。すべてのシーケンシングデータはhg19にマッピングされたので、これらの研究からVCFファイルをダウンロードし、正規のコーディング領域内でバリエーションを直接抽出した。

#### 【0582】

他の類人猿および哺乳類と比較するために、アカゲザル、マーモセット、ブタ、ウシ、ヤギ、ネズミ、およびニワトリを含む少数の他の種のSNPもdbSNPからダウンロードした。イヌ、ネコ、またはヒツジなどの他の種は廃棄した。それは、dbSNPがそれらの種に対して限られた数のバリエーションを提供するからである。最初に、各種のSNPをhg19にリフトオ

10

20

30

40

50

ーバーした。バリアントの約20%が偽遺伝子領域にマッピングされることが判明した。次いで、正規のコーディング領域の100種の脊椎動物の複数のアラインメントファイルから各種のエクソン座標を取得し、それらのエクソン内のバリアントを抽出した。次いで、それらの抽出されたSNPはhg19にリフトオーバーされた。バリアントがアラインメントとは異なる種のゲノムビルド上にある場合、まずアラインメントのゲノムビルドにバリアントをリフトした。

#### 【0583】

ウシSNPデータは様々な研究に由来するので、ウシバリアントのすべての大きなバッチをdbSNPからダウンロードし(VCFファイルが100MBより大きい16個のバッチ)、各バッチに対するミスセンス対同義比を計算することによってウシSNPの様々なバッチの品質を評価した。ミスセンス対同義比の中央値は0.781であり、中央絶対偏差(MAD)は0.160である(平均は0.879でありSDは0.496である)。異常値の比を伴う2つのバッチ(比が1.391であるsnpBatch\_1000\_BULL\_GENOMES\_1059190.gzおよび比が2.568であるsnpBatch\_COFACTOR\_GENOMICS\_1059634.gz)はさらなる分析から除外された。

10

#### 【0584】

[類人猿および哺乳類における多型の性質の評価]

大型類人猿SNPの有用性を実証するために、シングルTON SNPと一般的なSNP(アレル頻度(AF)>0.1%)の数の比を測定するエンリッチメントスコアを考案した。同義バリアントは、良性でありどのような選択圧力も受けずに一般に中立的に進化することが知られている。有害なミスセンスバリアントは、自然選択によって徐々に排除されるので、そのアレル頻度分布は同義バリアントと比較して稀なバリアントが多い傾向がある。

20

#### 【0585】

霊長類、哺乳類、および家禽において観察されるSNPと重複するgnomAD SNPに注目した。種毎の同義バリアントおよびミスセンスバリアントの数をカウントした。ミスセンスバリアントについては、「ミスセンス同一(missense identical)」と名付けられる、別の種において同一のアミノ酸変化を共有するタイプと、「ミスセンス相違(missense different)」と名付けられる、別の種において異なるアミノ酸変化を有するタイプという、2つのタイプへとさらに分類した。次いで、シングルTONバリアントの数と一般的なバリアントの数の比として、エンリッチメントスコアが種毎に計算された。

30

#### 【0586】

加えて、各種について同義バリアントとミスセンス同一バリアントとの間でエンリッチメントスコアを比較するために、2×2の分割表に対して相同性のカイ二乗( $\chi^2$ )検定を実行した。すべての霊長類が、同義バリアントとミスセンス同一バリアントとの間でエンリッチメントスコアに有意な差を示さないが、ウシ、ネズミ、およびニワトリは有意な差を示す。

#### 【0587】

この結果は、大型類人猿において同一のアミノ酸変化を共有するSNPが、同義SNPと非常に類似するエンリッチメントスコアを有する傾向があることを明らかにしており、それらのSNPに、ヒトの健康に対する軽度の影響があることを示唆している。異なるアミノ酸変化を有する、または大型類人猿において存在しないSNPは、同義SNPのエンリッチメントスコアから有意に逸脱するエンリッチメントスコアを有するが、非霊長類の種からのミスセンス多型も、同義バリアントと異なるアレル頻度分布を有する。結論は、大型類人猿において同一のアミノ酸変化を共有するSNPを、良性バリアントの訓練セットに追加することができるということである。

40

#### 【0588】

我々の仮定は、大半のバリアントが独立に派生したものであり、家系同一性(IBD)により生成されるのではないということである。したがって、IBD SNPにおける稀なバリアントのエンリッチメント分析を実行して、それらのエンリッチメントスコアの様々な挙動を評価した。IBD SNPは、チンパンジー、ボノボ、ゴリラ、S.オランウータン、およびB.オランウータンを含む、2つ以上の大型類人猿の種とヒトとの両方において現れる、ヒトSNP

50

として定義される。次いで、シングルTONの数を一般的なバリエーション(AF>0.1%)の数で割ったものとして定義されるエンリッチメントスコアが、ミスセンスバリエーションおよび同義バリエーションに対して別々に計算され、同義バリエーションは中立的であると考えられ比較の基準として働く。

【0589】

[ 哺乳類の種間での固定された置換 ]

固定された置換のエンリッチメント分析

種間の置換の稀なバリエーションエンリッチメント分析も研究した。UCSCゲノムブラウザ(<http://hgdownload.soe.ucsc.edu/goldenPath/hg19/multiz100way/hg19.100way.commonNames.nh>)から100種の脊椎動物の種の進化系統樹をダウンロードした。次いで、計算されたペア毎の進化系統的距離を計算し、近縁の種のペア(距離<0.3)を選択した。霊長類の種のペアを得るために、UCSCゲノムブラウザからCDS領域に対するヒトとの19種の哺乳類(16種の霊長類)ゲノムのアラインメント(hg38)をダウンロードした。4つの霊長類のペアが13個の脊椎動物のペアに追加された。以下の表は、一実装形態による、近縁の種の複数のペアの遺伝的距離を示す。

10

【0590】

【表 8】

種 1	種 2	距離
チンパンジー	ボノボ	0.00799195
アカゲザル	カニクイザル	0.00576807
マーモセット	リスザル	0.0680206
テングザル	キンシコウ	0.01339514
ヒヒ	ミドリザル	0.045652
ウマ	シロサイ	0.084397
ヒツジ	家畜のヤギ	0.1
マウス	ラット	0.176098
チンチラ	フサオネズミ	0.16
ケープキンモグラ	テンレック	0.265936
チャイニーズハムスター	ゴールデンハムスター	0.08
ホオヒゲコウモリ	ココウモリ	0.08254
イルカ	シャチ	0.074368
セーカーハヤブサ	ハヤブサ	0.2
ニワトリ	シチメンチョウ	0.126972
アオウミガメ	ニシキガメ	0.2
スッポン	トゲスッポン	0.2

10

20

30

40

## 【0591】

正規のコーディング領域内でヒトとの19種の哺乳類または99種の脊椎動物のゲノムの複数アラインメントを取り込み、脊椎動物の各々の選択されたペア間のヌクレオチド置換を得た。これらの置換は、種のペアとヒトバリエーションとの間のコドン変化が同一であることを条件として、gnomADからヒトエクソンSNPにマッピングされた。バリエーションを、同義バ

50

リアント、別の種において同一のアミノ酸変化を共有するミスセンスバリエーション、および別の種において異なるアミノ酸変化を有するミスセンスバリエーションという、3つのタイプへと分類した。エンリッチメントスコアが種のペア毎に各クラスに対して計算された。

#### 【0592】

##### 種内多型および種間多型の比較

チンパンジー、アカゲザル、マーモセット、ヤギ、ネズミ、およびニワトリを含む6つの種が、種内多型および種間多型の比較を実行するために選択され、それは、これらの種については種内バリエーションと種間バリエーションの両方が利用可能であったからである。種内バリエーションおよび種間バリエーションのエンリッチメントスコアの比較は、2つの2×2の分割表のオッズ比の比較に類似している。通常は、分割表間のオッズ比の相同性を評価するために、Woolf検定が適用される。したがって、Woolf検定を利用して、種内多型と種間多型との間のエンリッチメントスコアの差を評価した。

10

#### 【0593】

##### [ 遺伝子毎のエンリッチメント分析 ]

図64は、遺伝子毎のエンリッチメント分析の一実装形態を示す。一実装形態では、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器はさらに、病原性であると決定されたバリエーションの病原性を確認する遺伝子毎のエンリッチメント分析を実施するように構成される。遺伝的疾患を持つ個人の cohorts からサンプリングされた特定の遺伝子に対して、遺伝子毎のエンリッチメント分析は、病原性である特定の遺伝子におけるバリエーション候補を特定するために深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を適用することと、バリエーション候補の観察されるトリヌクレオチド変異率を合計してその合計を送信カウントおよび cohorts のサイズと乗じることに基づいて特定の遺伝子に対する変異の基準数を決定することと、病原性である特定の遺伝子の中の *de novo* ミスセンスバリエーションを特定するために深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を適用することと、変異の基準数を *de novo* ミスセンスバリエーションのカウントと比較することを含む。比較の出力に基づいて、遺伝子毎のエンリッチメント分析は、特定の遺伝子が遺伝子障害と関連付けられることと、*de novo* ミスセンスバリエーションが病原性であることとを確認する。いくつかの実装形態では、遺伝子障害は自閉スペクトラム障害(ASDと省略される)である。他の実装形態では、遺伝的障害は発達遅延障害(DDDと省略される)である。

20

30

#### 【0594】

図64に示される例では、特定の遺伝子の中の5つのバリエーション候補は、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器によって病原性であるものとして分類された。これらの5つのバリエーション候補は、 $10^{-8}$ 、 $10^{-2}$ 、 $10^{-1}$ 、 $10^5$ 、および $10^1$ という観察されたそれぞれのトリヌクレオチド変異率を有する。特定の遺伝子に対する変異の基準数は、5つのバリエーション候補のそれぞれの観察されたトリヌクレオチド変異率を合計し、その合計を送信/染色体カウント(2)および cohorts のサイズ(1000)と乗じることに基づいて、 $10^{-5}$ であると決定される。これが次いで *de novo* バリエーションカウント(3)と比較される。

#### 【0595】

いくつかの実装形態では、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器はさらに、出力としてp値を生み出す統計的検定を使用して比較を実行するように構成される。

40

#### 【0596】

他の実装形態では、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器はさらに、変異の基準数を *de novo* ミスセンスバリエーションのカウントと比較し、比較の出力に基づいて、特定の遺伝子が遺伝的疾患と関連付けられないことと、*de novo* ミスセンスバリエーションが良性であることとを確認するように構成される。

#### 【0597】

##### [ ゲノムワイドエンリッチメント分析 ]

50



図65は、ゲノムワイドエンリッチメント分析の一実装形態を示す。別の実装形態では、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器はさらに、病原性と決定されたバリエーションの病原性を確認するゲノムワイドエンリッチメント分析を実施するように構成される。ゲノムワイドエンリッチメント分析は、健康な個人の cohorts からサンプリングされた複数の遺伝子において病原性である de novo ミスセンスバリエーションの第1のセットを特定するために深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を適用することと、遺伝子障害を持つ個人の cohorts からサンプリングされる複数の遺伝子において病原性である de novo ミスセンスバリエーションの第2のセットを特定するために深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を適用することと、第1のセットおよび第2のセットのそれぞれのカウントを比較することと、比較の出力に基づいて、de novo ミスセンスバリエーションの第2のセットが遺伝的障害を持つ個人の cohorts においてエンリッチされ、したがって病原性であることを確認することを含む。いくつかの実装形態では、遺伝的疾患は自閉スペクトラム障害(ASDと省略される)である。他の実装形態では、遺伝的障害は発達遅延障害(DDDと省略される)である。

10

【0598】

いくつかの実装形態では、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器はさらに、p値を出力として生み出す統計的検定を使用して比較を実行するように構成される。一実装形態では、比較はさらにそれぞれの cohort サイズによってパラメータ化される。

20

【0599】

いくつかの実装形態では、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器はさらに、第1のセットおよび第2のセットのそれぞれのカウントを比較し、比較の出力に基づいて、de novo ミスセンスバリエーションの第2のセットが遺伝的障害を持つ個人の cohorts においてエンリッチされず、したがって良性であることを確認するように構成される。

【0600】

図65に示される例では、健康な cohorts における変異率(0.001)および影響を受けている cohorts における変異率(0.004)が、個人毎の変異率(4)とともに示されている。

【0601】

[具体的な実装形態]

バリエーション病原性分類器を構築するためのシステム、方法、および製造物品を説明する。実装形態の1つまたは複数の特徴は基本の実装形態と組み合わせられ得る。相互に排他的ではない実装形態は合成可能であると教示される。実装形態の1つまたは複数の特徴は他の実装形態と合成され得る。本開示は定期的にこれらの選択肢をユーザに思い起こさせる。これらの選択肢を繰り返し述べる記載がいくつかの実装形態において省略されていることは、先行するセクションにおいて教示された合成を限定するものと解釈されるべきではなく、これらの記載は以後の実装形態の各々へと前方に参照によって組み込まれる。

30

【0602】

開示される技術のシステム実装形態は、メモリに結合される1つまたは複数のプロセッサを含む。メモリは、ゲノム配列(たとえば、ヌクレオチド配列)におけるスプライスサイトを特定するスプライスサイト検出器を訓練するためのコンピュータ命令をロードされる。

40

【0603】

図48および図19に示されるように、システムは畳み込みニューラルネットワークベースのバリエーション病原性分類器を訓練し、これはメモリに結合される多数のプロセッサ上で実行される。システムは、良性バリエーションおよび病原性バリエーションから生成されたタンパク質配列ペアの、良性訓練例および病原性訓練例を使用する。良性バリエーションは、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準コドン配列を共有する代替のヒト以外の霊長類コドン配列上で発生するヒト以外の霊長類ミスセンスバリエーションを含む。「タンパク質配列ペア」という語句は、基準タンパク質配列および代替タンパク質配列を指し

50

、基準タンパク質配列は基準トリプレットヌクレオチド塩基(基準コドン)によって形成される基準アミノ酸を備え、代替タンパク質配列は代替トリプレットヌクレオチド塩基(代替コドン)によって形成される代替アミノ酸を備えるので、代替タンパク質配列は基準タンパク質配列の基準アミノ酸を形成する基準トリプレットヌクレオチド塩基(基準コドン)において発生するパリアントの結果として作り出される。パリアントは、SNP、挿入、または欠失であり得る。

【0604】

このシステムの実装形態および開示される他のシステムは任意選択で、以下の特徴のうちの一つまたは複数を含む。システムはまた、開示される方法に関連して説明される特徴を含み得る。簡潔にするために、システム特徴の代替的な組合せは個別に列挙されない。システム、方法、製造物品に適用可能な特徴は、基本の特徴の各々のstatutory classのセットに対して繰り返されない。このセクションにおいて特定される特徴がどのように他のstatutory classの中の基本の特徴と容易に組み合わせられ得るかを、読者は理解するであろう。

10

【0605】

図44に示されるように、一般的なヒトミスセンスパリアントは、少なくとも100000人のヒトからサンプリングされたヒト集団パリアントデータセットにわたって0.1%より高いマイナーアレル頻度(MAFと省略される)を有する。

【0606】

図44に示されるように、サンプリングされたヒトは異なるヒト亜集団に属し、一般的なヒトミスセンスパリアントはそれぞれのヒト亜集団パリアントデータセット内で0.1%より高いMAFを有する。

20

【0607】

ヒト亜集団は、アフリカ人/アフリカ系アメリカ人(AFRと省略される)、アメリカ人(AMRと省略される)、アシュケナーズ系ユダヤ人(ASJと省略される)、東アジア人(EASと省略される)、フィンランド人(FINと省略される)、フィンランド人以外のヨーロッパ人(NFEと省略される)、南アジア人(SASと省略される)、および他(OTHと省略される)を含む。

【0608】

図43および図44に示されるように、ヒト以外の霊長類ミスセンスパリアントは、チンパンジー、ボノボ、ゴリラ、B.オランウータン、S.オランウータン、アカゲザル、およびマーモセットを含む、複数のヒト以外の霊長類の種からのミスセンスパリアントを含む。

30

【0609】

図45および図46に示されるように、エンリッチメント分析に基づいて、システムは、ある特定のヒト以外の霊長類の種を、良性パリアントにその特定のヒト以外の霊長類の種のミスセンスパリアントを含めるために受け入れる。エンリッチメント分析は、特定のヒト以外の霊長類の種に対して、特定のヒト以外の霊長類の種の同義パリアントの第1のエンリッチメントスコアを、特定のヒト以外の霊長類の種のミスセンス同一パリアントの第2のエンリッチメントスコアと比較することを含む。

【0610】

図45は、ヒトオーソログミスセンスSNPの一実装形態を示す。ヒト以外の種におけるミスセンスSNPは、ヒトと一致する基準コドンおよび代替コドンを有する。図45に示されるように、ミスセンス同一パリアントは、ヒトと一致する基準コドン配列および代替コドン配列を共有するミスセンスパリアントである。

40

【0611】

図46および図47に示されるように、第1のエンリッチメントスコアは、MAFが0.1%より大きい一般的な同義パリアントに対する、MAFが0.1%より小さい稀な同義パリアントの比を決定することによって、作り出される。第2のエンリッチメントスコアは、MAFが0.1%より大きい一般的なミスセンス同一パリアントに対する、MAFが0.1%より小さい稀なミスセンス同一パリアントの比を決定することによって、作り出される。稀なパリアントは、シングルトンパリアントを含む。

50

## 【0612】

図46および図47に示されるように、第1のエンリッチメントスコアと第2のエンリッチメントスコアの差は所定の範囲内にあり、良性バリエーションに特定のヒト以外の霊長類のミスセンスバリエーションを含めるために、その特定のヒト以外の霊長類の種を受け入れることをさらに含む。差が所定の範囲にあることは、ミスセンス同一バリエーションが同義バリエーションと同じ程度の自然選択を受けているので、同義バリエーションと同じくらい良性であることを示す。

## 【0613】

図48に示されるように、システムは、良性バリエーションにヒト以外の霊長類の種のミスセンスバリエーションを含めるために、複数のヒト以外の霊長類の種を受け入れるようにエンリッチメント分析を繰り返し適用する。システムはさらに、ヒト以外の霊長類の種の各々に対する同義バリエーションの第1のエンリッチメントスコアとミスセンス同一バリエーションの第2のエンリッチメントスコアを比較するための、相同性のカイ二乗検定を含む。

10

## 【0614】

図48に示されるように、ヒト以外の霊長類ミスセンスバリエーションのカウントは少なくとも100000である。ヒト以外の霊長類ミスセンスバリエーションのカウントは385236である。一般的なヒトミスセンスバリエーションのカウントは少なくとも50000である。一般的なヒトミスセンスバリエーションのカウントは83546である。

## 【0615】

他の実装形態は、上で説明されたシステムの活動を実行するようにプロセッサによって実行可能な命令を記憶する、非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、上で説明されたシステムの活動を実行する方法を含み得る。

20

## 【0616】

開示される技術の別のシステムの実装形態は、一塩基多型(SNPと省略される)病原性分類器を構築することを含む。システムは畳み込みニューラルネットワークベースのSNP病原性分類器を訓練し、これは、良性SNPおよび病原性SNPによって表されるアミノ酸配列の良性訓練例および病原性訓練例を使用して、メモリに結合された多数のプロセッサ上で実行される。良性訓練例は、アミノ酸配列ペアとして表されるヌクレオチド配列の第1のセットおよび第2のセットを含み、各アミノ酸配列は、上流および下流のアミノ酸が側にある中心アミノ酸を含む。各アミノ酸配列ペアは、基準ヌクレオチド配列によって表されるアミノ酸の基準配列と、SNPを含む代替ヌクレオチド配列によって表されるアミノ酸の代替配列とを含む。

30

## 【0617】

図9に示されるように、第1のセットはヒトヌクレオチド配列ペアを備え、各ペアは、SNPを含みヒト集団内で一般的であると見なされるマイナーアレル頻度(MAFと省略される)を有するヒト代替ヌクレオチド配列を含む。第2のセットは、ヒト以外の霊長類代替ヌクレオチド配列とペアにされたヒト以外の霊長類基準ヌクレオチド配列を備える。ヒト以外の霊長類基準ヌクレオチド配列は、オーソロガスなヒトヌクレオチド基準配列を有する。ヒト以外の霊長類代替ヌクレオチド配列はSNPを含む。

## 【0618】

第1のシステムの実装形態についてこの特定の实装セクションにおいて論じられる特徴の各々は、このシステムの実装形態に等しく適用される。上で示されたように、すべてのシステム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

40

## 【0619】

他の実装形態は、上で説明されたシステムの活動を実行するようにプロセッサによって実行可能な命令を記憶する、非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、上で説明されたシステムの活動を実行する方法を含み得る。

## 【0620】

図48および図19に示されるように、開示される技術の第1の方法の実装形態は、バリア

50

ント病原性分類器を構築するステップを含む。方法はさらに、畳み込みニューラルネットワークベースのバリエーション病原性分類器を訓練するステップを含み、これは、良性バリエーションおよび病原性バリエーションから生成されるタンパク質配列ペアの良性訓練例および病原性訓練例を使用して、メモリに結合される多数のプロセッサ上で実行される。良性バリエーションは、一般的なヒトミスセンスバリエーションと、ヒトと一致する基準コドン配列を共有する代替的なヒト以外の霊長類コドン配列上で発生するヒト以外の霊長類ミスセンスバリエーションとを含む。

#### 【0621】

第1のシステム実装形態についてこの特定の実装セクションにおいて論じられる特徴の各々は、この方法の実装形態に等しく適用される。上で示されたように、すべてのシステム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

10

#### 【0622】

他の実装形態は、上で説明された方法を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、メモリと、上で説明された方法を実行するためにメモリに記憶された命令を実行するように動作可能な1つまたは複数のプロセッサとを含む、システムを含み得る。

#### 【0623】

図48および図19に示されるように、開示される技術の第2の方法の実装形態は、一塩基多型(SNPと省略される)病原性分類器を構築するステップを含む。方法はさらに、畳み込みニューラルネットワークベースのSNP病原性分類器を訓練するステップを含み、これは、良性SNPおよび病原性SNPによって表されるアミノ酸配列の良性訓練例および病原性訓練例を使用して、メモリに結合される多数のプロセッサ上で実行される。良性訓練例は、アミノ酸配列のペアとして表されるヌクレオチド配列の第1のセットおよび第2のセットを含み、各アミノ酸配列は上流および下流のアミノ酸が側にある中心アミノ酸を含み、各アミノ酸配列のペアは、基準ヌクレオチド配列によって表されるアミノ酸の基準配列およびSNPを含む代替ヌクレオチド配列によって表されるアミノ酸の代替配列を含む。第1のセットはヒトヌクレオチド配列ペアを備え、各ペアは、SNPを含みヒト集団内で一般的であると見なされるマイナーアレル頻度(MAFと省略される)を有する、ヒト代替ヌクレオチド配列を含む。第2のセットは、ヒト以外の霊長類代替ヌクレオチド配列とペアにされた、ヒト以外の霊長類基準ヌクレオチド配列を備える。ヒト以外の霊長類基準ヌクレオチド配列はオーソロガスなヒトヌクレオチド基準配列を有し、ヒト以外の霊長類代替ヌクレオチド配列はSNPを含む。

20

30

#### 【0624】

第2のシステム実装形態についてこの特定の実装セクションにおいて論じられる特徴の各々は、この方法の実装形態に等しく適用される。上で示されたように、すべてのシステム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

#### 【0625】

他の実装形態は、上で説明された方法を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、メモリと、上で説明された方法を実行するためにメモリに記憶された命令を実行するように動作可能な1つまたは複数のプロセッサとを含む、システムを含み得る。

40

#### 【0626】

二次構造および溶媒接触性分類器を伴う深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を使用するためのシステム、方法、および製造物品を説明する。実装形態の1つまたは複数の特徴は基本の実装形態と合成され得る。相互に排他的ではない実装形態は、合成可能であると教示される。実装形態の1つまたは複数の特徴は他の実装形態と合成され得る。本開示は定期的にこれらの選択肢をユーザに思い起こさせる。これらの選択肢を繰り返し述べる記載がいくつかの実装形態において省略されていることは、先行するセクションにおいて教示された合成を限定するものと解釈されるべきではなく、これらの記載は以後の実装形態の各々へと前方に参照によって組み込まれる。

50

## 【0627】

開示される技術のシステム実装形態は、メモリに結合される1つまたは複数のプロセッサを含む。メモリは、二次構造および溶媒接触性分類器を伴う深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を実行するためのコンピュータ命令をロードされる。

## 【0628】

システムは、タンパク質配列内のアミノ酸位置において3状態二次構造を予測するように訓練される、メモリに結合された多数のプロセッサ上で実行される第1の二次構造サブネットワークを備える。システムはさらに、タンパク質配列内のアミノ酸位置において3状態溶媒接触性を予測するように訓練される、メモリに結合された多数のプロセッサ上で実行される第2の溶媒接触性サブネットワークを備える。

10

## 【0629】

3状態二次構造は、複数のDNA二次構造状態であるヘリックス(H)、シート(B)、およびコイル(C)のうちの1つを指す。

## 【0630】

3状態溶媒接触性は、複数のタンパク質溶媒接触性である埋もれている(B)、中間(I)、および露出している(E)のうちの1つを指す。

## 【0631】

多数のプロセッサのうちの少なくとも1つで実行される位置特定の頻度行列(PFMと省略される)生成器は、霊長類PFM、哺乳類PFM、および脊椎動物PFMを生成するために、霊長類、哺乳類、および霊長類と哺乳類を除く脊椎動物の3つの配列グループに適用される。

20

## 【0632】

言い換えると、これは、PFM生成器を霊長類配列データに適用して霊長類PFMを生成すること、PFM生成器を哺乳類配列データに適用して哺乳類PFMを生成すること、PFM生成器を霊長類配列データおよび哺乳類配列データを含まない脊椎動物配列データに適用して脊椎動物PFMを生成することを含む。

## 【0633】

入力プロセッサは、標的バリエーションアミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸があるバリエーションアミノ酸配列を受け入れ、一塩基バリエーションが標的バリエーションアミノ酸を生み出す。多数のプロセッサのうちの少なくとも1つで実行される補足データ割振器は、そのバリエーションアミノ酸配列とアラインメントされた、標的基準アミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸がある基準アミノ酸配列を割り振る。これに続いて、補足データ割振器は、基準アミノ酸配列のために第1のサブネットワークおよび第2のサブネットワークによって作り出される基準状態分類を割り振る。この後で、補足データ割振器は、バリエーションアミノ酸配列のために第1のサブネットワークおよび第2のサブネットワークによって作り出されるバリエーション状態分類を割り振る。最後に、補足データ割振器は、基準アミノ酸配列とアラインメントされる霊長類PFM、哺乳類PFM、および脊椎動物PFMを割り振る。

30

## 【0634】

本出願の文脈では、「とアラインメントされる」という語句は、基準アミノ酸配列または代替アミノ酸配列の中の各アミノ場所に対する霊長類PFM、哺乳類PFM、および脊椎動物PFMを場所毎に決定し、アミノ酸場所が基準アミノ酸配列または代替アミノ酸配列において発生するのと同じ順序で場所毎にまたは場所の序数に基づいてその決定の結果を符号化して記憶することを指す。

40

## 【0635】

システムはまた、バリエーションアミノ酸配列、割り振られた基準アミノ酸配列、割り振られた基準状態分類およびバリエーション状態分類、ならびに割り振られたPFMを処理したことに基づいて、良性または病原性であるものとしてバリエーションアミノ酸配列を分類するように訓練された、多数のプロセッサで実行される深層畳み込みニューラルネットワークを含む。システムは、バリエーションアミノ酸配列に対する病原性スコアを少なくとも報告する出

50

カプロセッサを含む。

【0636】

このシステム実装形態および開示される他のシステムは任意選択で、以下の特徴のうちの1つまたは複数を含む。システムはまた、開示される方法に関連して説明される特徴を含み得る。簡潔にするために、システム特徴の代替的な組合せは個別に列挙されない。システム、方法、および製造物品に適用可能な特徴は、基本の特徴の各statutory classセットに対して繰り返されない。読者は、このセクションにおいて特定される特徴が他のstatutory classにおける基本の特徴とどのように容易に合成され得るかを理解するのである。

【0637】

深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を備えるシステムはさらに、病原性スコアに基づいて良性または病原性として一塩基バリエーションを分類するように構成される。

【0638】

システムは、深層畳み込みニューラルネットワークが、少なくとも、バリエーションアミノ酸配列、割り振られた基準アミノ酸配列、割り振られたバリエーション二次構造状態分類、割り振られた基準二次構造状態分類、割り振られたバリエーション溶媒接触性状態分類、割り振られた基準溶媒接触性状態分類、割り振られた霊長類PFM、割り振られた哺乳類PFM、および割り振られた脊椎動物PFMを、入力として並列に受け入れるような、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を備える。

【0639】

システムは、バッチ正規化層、ReLU非線形性層、および次元変更層を使用して、バリエーションアミノ酸配列、割り振られた基準アミノ酸配列、割り振られた霊長類PFM、割り振られた哺乳類PFM、および割り振られた脊椎動物PFMを前処理するように構成される。システムはさらに、前処理された特性評価を合計し、その合計を、割り振られたバリエーション二次構造状態分類、割り振られた基準二次構造状態分類、割り振られたバリエーション溶媒接触性状態分類、および割り振られた基準溶媒接触性状態分類と連結して、連結された入力を生み出すように構成される。システムは、次元変更層を通じて連結された入力を処理し、処理された連結された入力を受け入れて深層畳み込みニューラルネットワークの残差ブロックを開始する。

【0640】

深層畳み込みニューラルネットワークは、配列において最低から最高まで並べられる残差ブロックのグループを備える。深層畳み込みニューラルネットワークは、残差ブロックの数、スキップ接続の数、および非線形活性化を伴わない残差接続の数によってパラメータ化される。深層畳み込みニューラルネットワークは、先行する入力の空間次元および特徴量次元を形状変更する次元変更層を備える。

【0641】

システムはさらに、霊長類、哺乳類、および脊椎動物にわたってアライメントされた基準アミノ酸配列において保存されている標的基準アミノ酸から標的バリエーションアミノ酸を作り出す、一塩基バリエーションを病原性として分類するように訓練するように構成される。

【0642】

保存率は、標的基準アミノ酸の機能的な有意性を表し、PFWから決定される。システムはさらに、バリエーションアミノ酸配列と基準バリエーションアミノ酸配列との間で異なる二次構造を引き起こす一塩基バリエーションを病原性として分類するように訓練するように構成される。

【0643】

システムはさらに、バリエーションアミノ酸配列と基準バリエーションアミノ酸配列との間で異なる溶媒接触性を引き起こす一塩基バリエーションを病原性として分類するように訓練するように構成される。

【0644】

10

20

30

40

50

PFMは、他の種のアライメントされるタンパク質配列にわたるヒトタンパク質配列におけるアミノ酸の出現の頻度を位置ごとに決定することによって、他の種のアライメントされるタンパク質配列にわたるヒトタンパク質配列におけるアミノ酸の保存率を表す。

【0645】

二次構造の3状態は、ヘリックス、シート、およびコイルである。第1の二次構造サブネットワークは、入力タンパク質配列と、入力タンパク質配列内のアミノ酸位置とアライメントされる霊長類PFM、哺乳類PFM、および脊椎動物PFMとを受け入れ、アミノ酸位置の各々において3状態二次構造を予測するように訓練される。溶媒接触性の3状態は、露出している、埋もれている、および中間である。

【0646】

二次溶媒接触性サブネットワークは、入力タンパク質配列と、入力タンパク質配列内のアミノ酸位置とアラインメントされている霊長類PFM、哺乳類PFM、および脊椎動物PFMとを受け入れ、アミノ酸位置の各々において3状態溶媒接触性を予測するように訓練される。入力タンパク質配列は基準タンパク質配列である。入力タンパク質配列は代替タンパク質配列である。第1の二次構造サブネットワークは、配列において最低から最高まで並べられる残差ブロックのグループを備える。第1の二次構造サブネットワークは、残差ブロックの数、スキップ接続の数、および非線形活性化を伴わない残差接続の数によってパラメータ化される。

【0647】

第1の二次構造サブネットワークは、先行する入力空間次元および特徴量次元を形状変更する次元変更層を備える。第2の溶媒接触性サブネットワークは、配列において最低から最高まで並べられる残差ブロックのグループを備える。第2の溶媒接触性サブネットワークは、残差ブロックの数、スキップ接続の数、および非線形活性化を伴わない残差接続の数によってパラメータ化される。第2の溶媒接触性サブネットワークは、先行する入力空間次元および特徴量次元を形状変更する次元変更層を備える。

【0648】

各残差ブロックは、少なくとも1つのバッチ正規化層、少なくとも1つの正規化線形ユニット(ReLUと省略される)層、少なくとも1つの次元変更層、および少なくとも1つの残差接続を備える。各残差ブロックは、2つのバッチ正規化層、2つのReLU非線形性層、2つの次元変更層、および1つの残差接続を備える。

【0649】

深層畳み込みニューラルネットワーク、第1の二次構造サブネットワーク、および第2の溶媒接触性サブネットワークは各々、最終分類層を備える。最終分類層はシグモイドベースの層である。最終分類層はソフトマックスベースの層である。

【0650】

システムはさらに、深層畳み込みニューラルネットワークとの協調のために、第1の二次構造サブネットワークおよび第2の溶媒接触性サブネットワークの最終分類層を除去するように構成される。

【0651】

システムはさらに、深層畳み込みニューラルネットワークの訓練の間に、誤差をサブネットワークに逆伝播してサブネットワーク重みを更新することを含めて、第1の二次構造サブネットワークおよび第2の溶媒接触性サブネットワークを病原性分類についてさらに訓練するように構成される。

【0652】

第2の溶媒接触性サブネットワークは少なくとも膨張畳み込み層を備える。システムはさらに、発達遅延障害(DDDと省略される)を引き起こすバリエーションを病原性として分類するように構成される。バリエーションアミノ酸配列および基準アミノ酸配列はランキングアミノ酸を共有する。システムはさらに、深層畳み込みニューラルネットワークへの入力を符号化するためにワンホット符号化を使用するように構成される。

【0653】

10

20

30

40

50

図1Qは、開示される技術が動作することのできる例示的なコンピューティング環境を示す。深層畳み込みニューラルネットワーク、第1の二次構造サブネットワーク、および第2の溶媒接触性サブネットワークは、1つまたは複数の訓練サーバ上で訓練される。訓練された深層畳み込みニューラルネットワーク、第1の訓練された二次構造サブネットワーク、および訓練された第2の溶媒接触性サブネットワークは、要求側のクライアントから入力配列を受け取る1つまたは複数の本番サーバ上に展開される。本番サーバは、深層畳み込みニューラルネットワーク、第1の二次構造サブネットワーク、および第2の溶媒接触性サブネットワークのうちの少なくとも1つを通じて入力配列を処理して、クライアントに送信される出力を作り出す。

【0654】

他の実装形態は、上で説明されたシステムの活動を実行するようにプロセッサによって実行可能な命令を記憶する、非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、上で説明されたシステムの活動を実行する方法を含み得る。

【0655】

開示される技術の別のシステム実装形態は、メモリに結合された多数のプロセッサ上で実行される、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を含む。システムは、霊長類PFMおよび哺乳類PFMを生成するために霊長類および哺乳類の2つの配列グループに適用される、多数のプロセッサのうちの少なくとも1つで実行される、位置特定の頻度行列(PFMと省略される)生成器を含む。システムはまた、標的バリエーションアミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸があるバリエーションアミノ酸配列を受け入れる入力プロセッサを含み、一塩基バリエーションが標的バリエーションアミノ酸を生み出す。システムはまた、バリエーションアミノ酸配列とアラインメントされる、標的基準アミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸がある基準アミノ酸配列を割り振る、多数のプロセッサのうちの少なくとも1つで実行される補足データ割振器を含む。補足データ割振器はまた、基準アミノ酸配列とアラインメントされた霊長類PFMおよび哺乳類PFMを割り振る。システムはさらに、バリエーションアミノ酸配列、割り振られた基準アミノ酸配列、および割り振られたPFMを処理することに基づいてバリエーションアミノ酸配列を良性または病原性として分類するように訓練される、多数のプロセッサ上で実行される深層畳み込みニューラルネットワークを含む。最後に、システムは、バリエーションアミノ酸配列に対する病原性スコアを少なくとも報告する出力プロセッサを含む。

【0656】

このシステム実装形態および開示される他のシステムは任意選択で、以下の特徴のうちの1つまたは複数を含む。システムはまた、開示される方法に関連して説明される特徴を含み得る。簡潔にするために、システム特徴の代替的な組合せは個別に列挙されない。システム、方法、および製造物品に適用可能な特徴は、基本の特徴の各statutory classに対して繰り返されない。このセクションにおいて特定される特徴が他のstatutory classの中の基本の特徴とどのように容易に組み合わせられ得るかを、読者は理解するであろう。

【0657】

システムはさらに、病原性スコアに基づいて一塩基バリエーションを良性または病原性として分類するように構成される。深層畳み込みニューラルネットワークは、バリエーションアミノ酸配列、割り振られた基準アミノ酸配列、割り振られた霊長類PFM、および割り振られた哺乳類PFMを並列に受け入れて処理する。システムはさらに、霊長類および哺乳類にわたって基準アミノ酸配列において保存されている標的基準アミノ酸から標的バリエーションアミノ酸を作り出す、一塩基バリエーションを病原性として分類するように訓練するように構成される。保存率は、標的基準アミノ酸の機能的な有意性を表し、PFWから決定される。

【0658】

第1のシステム実装形態に対してこの特定の実装セクションにおいて論じられる特徴の各々はこのシステム実装形態に等しく適用される。上で示されたように、すべてのシステム特徴は、ここでは繰り返されず、参照によって繰り返されるものと見なされるべきであ

10

20

30

40

50



る。

【0659】

他の実装形態は、上で説明されたシステムの活動を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、上で説明されたシステムの活動を実行する方法を含み得る。

【0660】

開示される技術の第1の方法の実装形態は、タンパク質配列内のアミノ酸位置において3状態二次構造を予測するように訓練される、メモリに結合された多数のプロセッサ上で第1の二次構造サブネットワークを実行するステップを含む。タンパク質配列内のアミノ酸位置において3状態溶媒接触性を予測するように訓練される、メモリに結合される多数のプロセッサ上で第2の溶媒接触性サブネットワークを実行すること。多数のプロセッサのうち少なくとも1つで実行すること。霊長類位置特定の頻度行列(PFMと省略される)、哺乳類PFM、および脊椎動物PFMを生成するために、霊長類、哺乳類、および霊長類と哺乳類を除く脊椎動物の3つの配列グループに適用される、PFM生成器。標的バリエーションアミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸があるバリエーションアミノ酸配列入力プロセッサを受け入れること。一塩基バリエーションは標的バリエーションアミノ酸を作り出す。バリエーションアミノ酸配列とアラインメントされた、標的基準アミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸がある基準アミノ酸配列を割り振る補足データ割振器を、多数のプロセッサのうち1つで実行すること。補足データ割振器はまた、基準アミノ酸配列のために第1のサブネットワークおよび第2のサブネットワークによって作り出される基準状態分類を割り振る。補足データ割振器はさらに、バリエーションアミノ酸配列のために第1のサブネットワークおよび第2のサブネットワークによって作り出されるバリエーション状態分類を割り振る。補足データ割振器は、基準アミノ酸配列とアラインメントされた、霊長類PFM、哺乳類PFM、および脊椎動物PFMを割り振る。バリエーションアミノ酸配列、割り振られた基準アミノ酸配列、割り振られた基準状態分類およびバリエーション状態分類、ならびに割り振られたPFMを処理することに基づいて、バリエーションアミノ酸配列を良性または病原性として分類するように訓練される深層畳み込みニューラルネットワークを、多数のプロセッサ上で実行すること。出力プロセッサを通じてバリエーションアミノ酸配列に対する病原性スコアを少なくとも報告すること。

【0661】

第1のシステム実装形態に対するこの特定の実装セクションにおいて論じられる特徴の各々はこの方法の実装形態に等しく適用される。上で示されたように、すべてのシステム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

【0662】

他の実装形態は、上で説明された方法を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、メモリと、上で説明された方法を実行するためにメモリに記憶された命令を実行するように動作可能な1つまたは複数のプロセッサとを含む、システムを含み得る。

【0663】

開示される技術の第2の方法の実装形態は、深層畳み込みニューラルネットワークベースのバリエーション病原性分類器を、メモリに結合された多数のプロセッサ上で実行するステップを含む。霊長類PFMおよび哺乳類PFMを生成するために霊長類と哺乳類の2つの配列グループに適用される、多数のプロセッサのうち少なくとも1つで位置特定の頻度行列(PFMと省略される)生成器を実行すること。入力プロセッサにおいて、標的バリエーションアミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸があるバリエーションアミノ酸配列を受け入れること。一塩基バリエーションは標的バリエーションアミノ酸を作り出す。バリエーションアミノ酸配列とアラインメントされる、標的基準アミノ酸の上流および下流の側に各方向への少なくとも25個のアミノ酸がある基準アミノ酸配列を割り振り、基準アミノ酸配列とアラインメントされる霊長類PFMおよび哺乳類PFMを割り振る、多数のプロセッサのうち少なくとも1つで補足データ割振器を実行すること。バリエーションアミノ酸配列、割

10

20

30

40

50

り振られた基準アミノ酸配列、および割り振られたPFMを処理することに基づいて、バリエーションアミノ酸配列を良性または病原性として分類するように訓練される深層畳み込みニューラルネットワークを、多数のプロセッサ上で実行すること。出力プロセッサにおいてバリエーションアミノ酸配列に対する病原性スコアを少なくとも報告すること。

【0664】

第2のシステム実装形態に対するこの特定の实装セクションにおいて論じられる特徴の各々は、この方法の実装形態に等しく適用される。上で示されたように、すべてのシステム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

【0665】

他の実装形態は、上で説明された方法を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、メモリと、上で説明された方法を実行するためにメモリに記憶された命令を実行するように動作可能な1つまたは複数のプロセッサとを含む、システムを含み得る。

【0666】

開示される技術のさらに別のシステム実装形態は、一塩基多型(SNPと省略される)病原性分類器を訓練するための大規模な病原性訓練データを生成するシステムを含む。

【0667】

図19に示されるように、システムは、良性SNPの訓練セットと、組合せで生成されたSNPの合成セットからカリング(cull)される予測されるエリート病原性SNPの訓練セットとを使用して、メモリに結合された多数のプロセッサ上で実行されるSNP病原性分類器を訓練する。本出願の文脈では、予測されるエリート病原性SNPは、アンサンプルによって出力されるような、平均病原性スコアまたは最大病原性スコアに基づいて各サイクルの終わりに作成/選択されるSNPである。「エリート」という用語は、遺伝的アルゴリズムの語彙から借用され、遺伝的アルゴリズムの出版物において通常与えられる意味を有することが意図される。

【0668】

図37、図38、図39、図40、図41、および図42に示されるように、システムはサイクルの中で反復的にエリートセットを構築する。このとき、予測されるSNPがない状態から始めて、合成セットから異常値SNPをカリングすることによって予測されるSNPの完全なセットを累積する。合成セットは、良性セットに存在しない、組合せで生成されるSNPである疑似病原性SNPを備え、異常値SNPがエリートセットへの包含のために合成セットから反復的にカリングされるにつれてセットのメンバー数が減少する。本出願の文脈では、「カリングする」という用語は、以前の集団をフィルタリングすること、新しい集団で置き換えること、更新すること、または選択することを意味する。「カリングする」という用語は、遺伝的アルゴリズムの語彙から借用され、遺伝的アルゴリズムの出版物において通常与えられる意味を有することが意図される。

【0669】

図37、図38、図39、図40、図41、および図42に示されるように、システムは、サイクルにおいて反復的に、合成セットから異常値SNPをカリングするために、SNP病原性分類器のアンサンプルを訓練して適用する。これは、良性SNPの一般訓練セット、予測されるエリート病原性SNPの一般訓練セット、および置換を伴わずに合成セットからサンプリングされる疑似病原性SNPの別個の訓練セットを使用して、アンサンプルを訓練することを含む。これはまた、現在のサイクルにおいてアンサンプルを訓練するために使用されなかった合成セットからの少なくともいくつかのSNPをスコアリングするために、訓練されたアンサンプルを適用し、スコアリングされたSNPから、一般エリートセットにおいて累積すべき現在のサイクルの異常値SNPをスコアリングされたSNPから選択するためにスコアを使用することによって、合成セットから異常値SNPをカリングしてカリングされた異常値SNPを一般エリートセットにおいて累積するように、訓練されたアンサンプルを適用することを含む。

【0670】

10

20

30

40

50

本出願の文脈では、「疑似病原性SNP」は、訓練の目的で病原性としてラベリングされ、訓練の間に置換を伴わずに合成的に生成されたバリエーションからサンプリングされるSNPである。

【0671】

また、予測されるエリート病原性SNPの訓練セットは、複数のサイクルにわたって反復的に構築される。

【0672】

図37、図38、図39、図40、図41、および図42に示されるように、システムは次いで、訓練によって導かれた分類器パラメータ、複数のサイクルにわたって完成され一般良性子集の所定の範囲内にある一般エリートセットと、SNP病原性分類器を訓練するための一般良性子集とを、メモリに記憶する。

10

【0673】

図37、図38、図39、図40、図41、および図42に示されるように、予測されるエリート病原性SNPは、アンサンプルによって予測されるSNPの上位5%である。いくつかの実装形態では、それらは20000個などの固定された数の上位にスコアリングされるSNPである。

【0674】

SNP病原性分類器およびSNP病原性分類器のアンサンプルは各々、深層畳み込みニューラルネットワーク(DCNNと省略される)である。アンサンプルは4個から16個のDCNNを含む。図37、図38、図39、図40、図41、および図42に示されるように、アンサンプルは8個のDCNNを含む。

20

【0675】

図37、図38、図39、図40、図41、および図42に示されるように、システムは、サイクルの間のエポックにおいてDCNNのアンサンプルを訓練し、妥当性確認サンプルに対する予測が良性予測と病原性予測の別々の確率分布クラスタを形成するとき、特定のサイクルに対する訓練を終了する。

【0676】

図37、図38、図39、図40、図41、および図42に示されるように、システムは、スコアを使用して、DCNNのアンサンプルからのスコアを合計することによって現在のサイクルの異常値SNPを選択する。

【0677】

図37、図38、図39、図40、図41、および図42に示されるように、システムは、スコアを使用して、DCNNのアンサンプルによってスコアリングされるSNPの各々に対する最大平均値をとることによって現在のサイクルの異常値SNPを選択する。

30

【0678】

図37、図38、図39、図40、図41、および図42に示されるように、現在のサイクルの間の置換を伴わないサンプリングは、現在のサイクルの間の疑似病原性SNPの互いに素の別々の訓練セットをもたらす。

【0679】

システムは、終了条件に達するまでサイクルを続ける。終了条件はサイクルの所定の数であり得る。図37、図38、図39、図40、図41、および図42に示されるように、サイクルの所定の数は21である。

40

【0680】

図37、図38、図39、図40、図41、および図42に示されるように、終了条件は、予測されるエリート病原性セットサイズが良性セットサイズの所定の範囲内にあるときである。

【0681】

分類器パラメータは、少なくとも畳み込みフィルタ重みおよび学習率であり得る。

【0682】

システムは、SNP病原性分類器としてアンサンプルの中のSNP病原性分類器のうちの1つを選択することができる。選択されるSNP病原性分類器は、最終のサイクルにおいて評価される妥当性確認サンプルについてアンサンプルの中の他のSNP病原性分類器より予測が

50

優れていた分類器であり得る。

【0683】

図37、図38、図39、図40、図41、および図42に示されるように、複数のサイクルにわたって完成する一般エリートセットは、少なくとも400000個の予測されるエリート病原性SNPを有し得る。

【0684】

図37、図38、図39、図40、図41、および図42に示されるように、システムは、各サイクルにおいて、予測されるエリート病原性SNPにおける変異率バイアスを防ぐために、良性SNPとサンプリングされた疑似病原性SNPとの間でトリヌクレオチドコンテキストを照合することができる。

10

【0685】

図37、図38、図39、図40、図41、および図42に示されるように、同期セットからの疑似病原性SNPのサンプリングは、各々の連続するサイクルにおいて5%ずつ減少し得る。

【0686】

図37、図38、図39、図40、図41、および図42に示されるように、システムは、訓練のために現在のサイクルにおいてサンプリングされる疑似病原性SNPによって現在のサイクルにおいてスコアリングされる合成SNP、予測されるエリート病原性SNP、および訓練のために現在のサイクルにおいて使用される良性SNPをフィルタリングすることができる。

【0687】

第1のシステム実装形態に対するこの特定の实装セクションにおいて論じられる特徴の各々が、このシステム実装形態に等しく適用される。上で示されるように、すべてのシステム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

20

【0688】

他の実装形態は、上で説明されたシステムの活動を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、メモリと、上で説明されたシステムの活動を実行するようにメモリに記憶されている命令を実行するように動作可能な1つまたは複数のプロセッサとを含む、システムを含み得る。

【0689】

開示される技術の別の実装形態は、図36に示されるように、畳み込みニューラルネットワーク(CNNと省略される)ベースの半教師あり学習器を含む。

30

【0690】

図36に示されるように、半教師あり学習器は、良性訓練セットおよび病原性訓練セットについて反復的に訓練される、メモリに結合された多数のプロセッサ上で実行されるCNNのアンサンブルを含み得る。

【0691】

図36に示されるように、半教師あり学習器は、訓練されたアンサンブルによる合成セットの評価に基づいて病原性訓練セットのセットサイズを次第に増強する、プロセッサのうちの少なくとも1つで実行されるセット増強器を含み得る。

40

【0692】

各反復において、評価は、セット増強器によって病原性訓練セットに追加される、予測されるエリート病原性セットを作り出す。

【0693】

半教師あり学習器は、CNN、増強された病原性訓練セット、および良性訓練セットのうちの少なくとも1つを使用して、一塩基多型(SNPと省略される)病原性分類器を構築して訓練する、ビルダを含み得る。

【0694】

第1のシステム実装形態に対するこの特定の实装セクションにおいて論じられる特徴の各々が、このシステム実装形態に等しく適用される。上で示されるように、すべてのシス

50

テム特徴はここで繰り返されず、参照によって繰り返されるものと見なされるべきである。

【0695】

他の実装形態は、上で説明されたシステムの活動を実行するようにプロセッサによって実行可能な命令を記憶する非一時的コンピュータ可読記憶媒体を含み得る。さらに別の実装形態は、メモリと、上で説明されたシステムの活動を実行するようにメモリに記憶されている命令を実行するように動作可能な1つまたは複数のプロセッサとを含む、システムを含み得る。

【0696】

先行する説明は、開示される技術の作成および使用を可能にするために提示される。開示される実装形態に対する様々な修正が明らかであり、本明細書で定義される一般原理は、開示される技術の趣旨および範囲から逸脱することなく、他の実装形態および適用例に適用され得る。したがって、開示される技術は、示される実装形態に限定されることは意図されず、本明細書で開示される原理および特徴と一致する最も広い範囲を認められるべきである。開示される技術の範囲は添付の特許請求の範囲によって定義される。

【0697】

[コンピュータシステム]

図66は、開示される技術を実装するために使用され得るコンピュータシステムの簡略化されたブロック図である。コンピュータシステムは通常、バスサブシステムを介していくつかの周辺デバイスと通信する少なくとも1つのプロセッサを含む。これらの周辺デバイスは、たとえば、メモリデバイスおよびファイルストレージサブシステム、ユーザインターフェース入力デバイス、ユーザインターフェース出力デバイス、ならびにネットワークインターフェースサブシステムを含む、ストレージサブシステムを含み得る。入力デバイスおよび出力デバイスはコンピュータシステムとのユーザの対話を可能にする。ネットワークインターフェースサブシステムは、他のコンピュータシステムにおける対応するインターフェースデバイスへのインターフェースを含む、外部ネットワークへのインターフェースを提供する。

【0698】

一実装形態では、良性データセット生成器、バリエーション病原性分類器、二次構造分類器、溶媒接触性分類器、および半教師あり学習器などのニューラルネットワークは、ストレージサブシステムおよびユーザインターフェース入力デバイスへ通信可能につながる。

【0699】

ユーザインターフェース入力デバイスは、キーボードと、マウス、トラックボール、タッチパッド、またはグラフィックタブレットなどのポインティングデバイスと、ディスプレイに組み込まれたタッチスクリーンと、音声認識システムおよびマイクロフォンなどのオーディオ入力デバイスと、他のタイプの入力デバイスとを含み得る。一般に、「入力デバイス」という用語の使用は、コンピュータシステムへ情報を入力するためのすべての可能なタイプのデバイスおよび方式を含むことが意図される。

【0700】

ユーザインターフェース出力デバイスは、ディスプレイサブシステム、プリンタ、faxマシン、またはオーディオ出力デバイスなどの非視覚的ディスプレイを含み得る。ディスプレイサブシステムは、陰極線管(CRT)、液晶ディスプレイ(LCD)などのフラットパネルデバイス、プロジェクションデバイス、または可視の画像を創造するための何らかの他の機構を含み得る。ディスプレイサブシステムはまた、オーディオ出力デバイスなどの非視覚ディスプレイを提供することができる。一般に、「出力デバイス」という用語の使用は、コンピュータシステムから情報をユーザまたは別の機械もしくはコンピュータシステムに出力するためのすべての可能なタイプのデバイスおよび方式を含むことが意図される。

【0701】

ストレージサブシステムは、本明細書で説明されるモジュールおよび方法の一部またはすべての機能を提供する、プログラミングおよびデータ構築物を記憶する。これらのソフ

10

20

30

40

50

トウェアモジュールは一般に、プロセッサだけによって、または他のプロセッサと組み合わせて実行される。

【0702】

ストレージサブシステムにおいて使用されるメモリは、プログラム実行の間の命令およびデータの記憶のためのメインランダムアクセスメモリ(RAM)と、固定された命令が記憶される読取り専用メモリ(ROM)とを含む、いくつかのメモリを含み得る。ファイルストレージサブシステムは、プログラムおよびデータファイルのための永続的なストレージを提供することができ、ハードディスクドライブ、関連する取り外し可能なメディアを伴うフロッピーディスクドライブ、CD-ROMドライブ、光学ドライブ、または取り外し可能なメディアカートリッジを含み得る。いくつかの実装形態の機能を実装するモジュールは、ストレージサブシステムの中の、または他のプロセッサによってアクセス可能な他の機械の中の、ファイルストレージサブシステムによって記憶され得る。

10

【0703】

バスサブシステムは、コンピュータシステムの様々な構成要素およびサブシステムに意図されるように互いに通信させるための機構を提供する。バスサブシステムは単一のバスとして概略的に示されているが、バスサブシステムの代替的な実装形態は複数のバスを使用することができる。

【0704】

コンピュータシステム自体が、パーソナルコンピュータ、ポータブルコンピュータ、ワークステーション、コンピュータ端末、ネットワークコンピュータ、テレビジョン、メインフレーム、サーバファーム、緩やかにネットワーク化されたコンピュータの広く分布するセット、または、任意の他のデータ処理システムもしくはユーザデバイスを含む、様々なタイプであってよい。コンピュータおよびネットワークの変わり続ける性質により、図66に示されるコンピュータシステムの記述は、開示される技術を例示することを目的とする特定の例としてのみ意図されている。図66に示されるコンピュータシステムより多数または少数の構成要素を有する、コンピュータシステムの多くの他の構成が可能である。

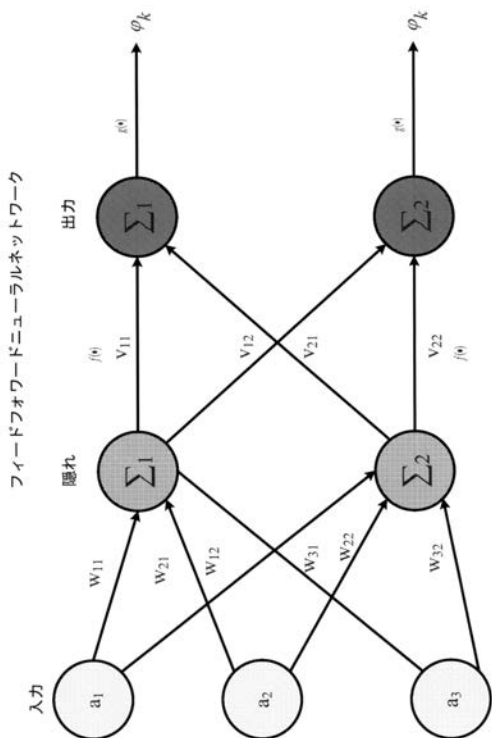
20

【0705】

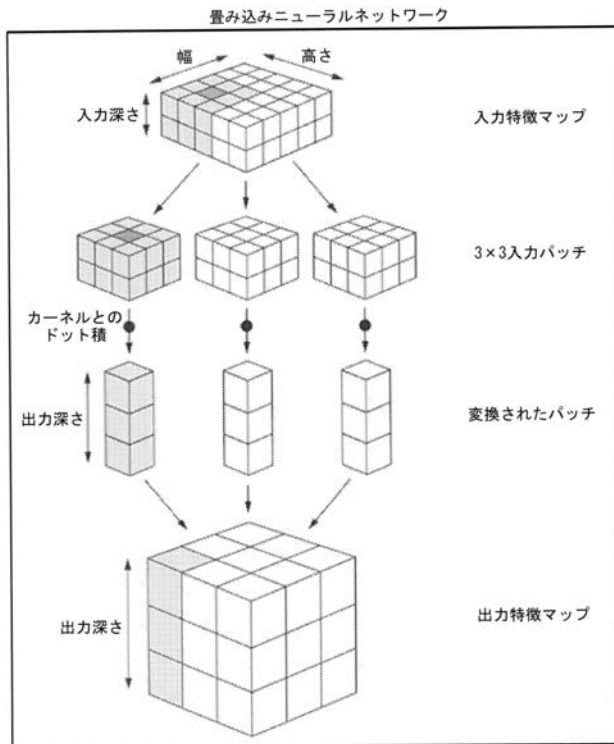
深層学習プロセッサは、GPUまたはFPGAであってよく、Google Cloud Platform、Xilinx、およびCirrascalなどの深層学習クラウドプラットフォームによってホストされてよい。深層学習プロセッサの例には、GoogleのTensor Processing Unit(TPU)、GX4 Rackmount Series、GX8 Rackmount Seriesのようなラックマウントソリューション、NVIDIA DGX-1、MicrosoftのStratix V FPGA、GraphcoreのIntelligent Processor Unit(IPU)、Snapdragonプロセッサを用いたQualcommのZerothプラットフォーム、NVIDIAのVolta、NVIDIAのDRIVE PX、NVIDIAのJETSON TX1/TX2 MODULE、IntelのNirvana、Movidius VPU、Fujitsu DPI、ARMのDynamicIQ、IBM TrueNorthなどがある。

30

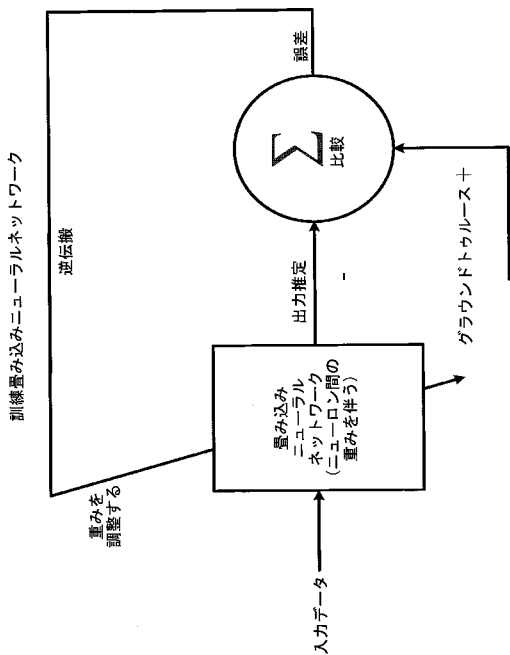
【図 1 A】



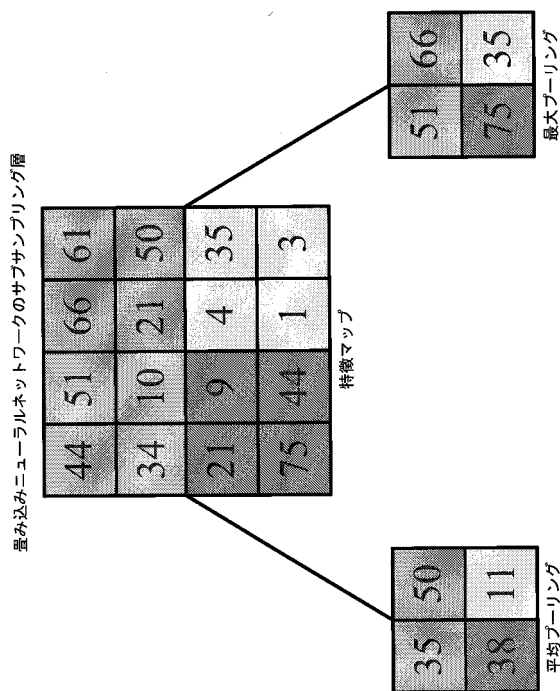
【図 1 B】



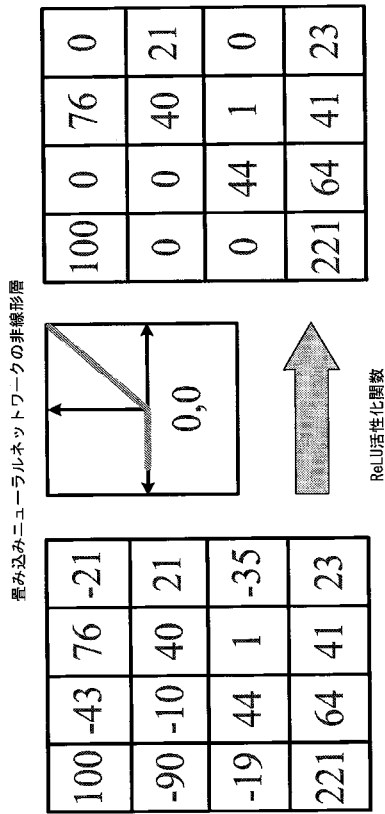
【図 1 C】



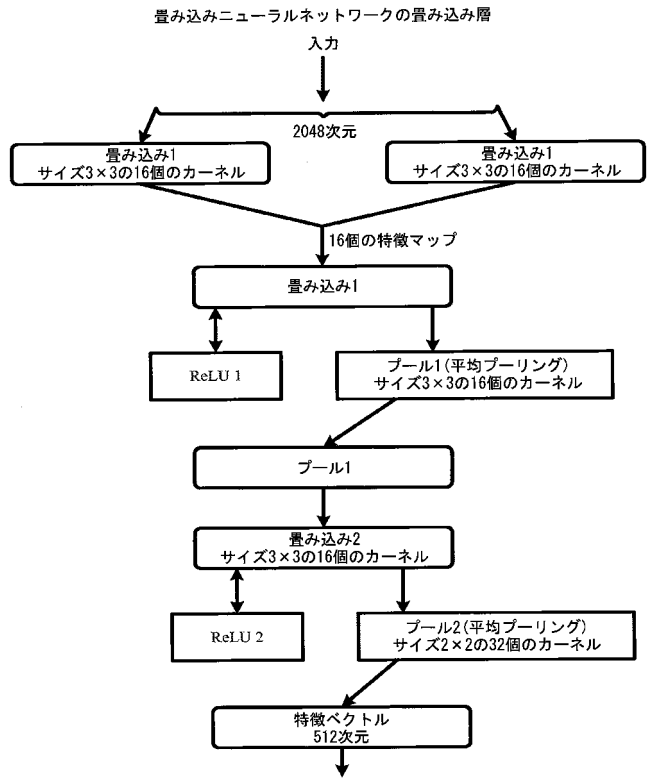
【図 1 D】



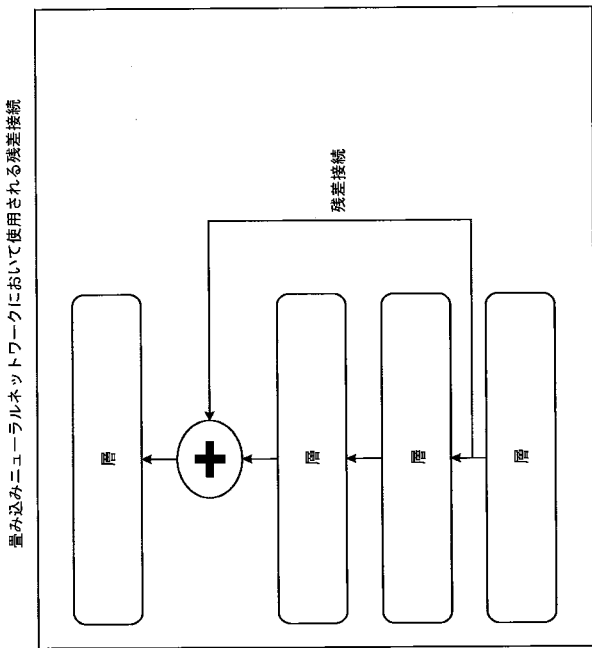
【図 1 E】



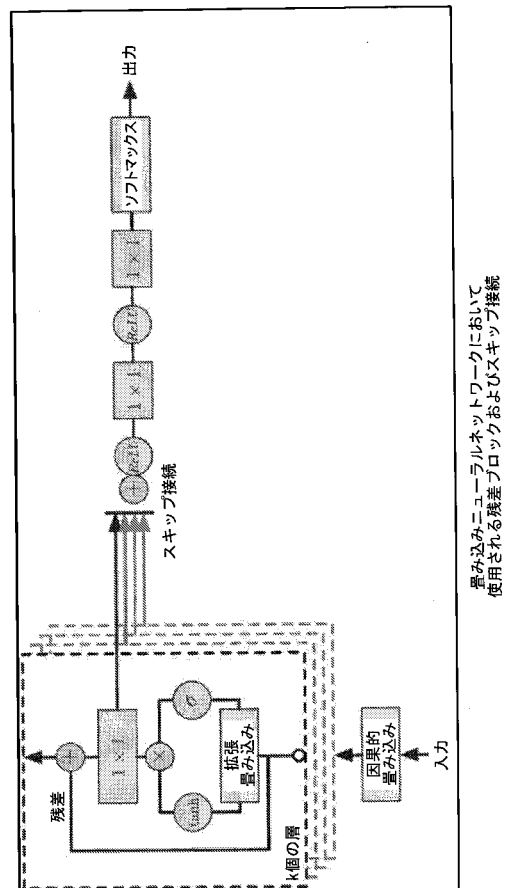
【図 1 F】



【図 1 G】



【図 1 H】





【 図 1 I 】

畳み込みニューラルネットワークを用いたバッチ正規化フォワードパス

$$\mu_B = \frac{1}{n} \sum_{i=1}^n x_i^{(\ell-1)}$$

$$\sigma_B^2 = \frac{1}{n} \sum_{i=1}^n (x_i^{(\ell-1)} - \mu_B)^2$$

$$\hat{x}^{(\ell-1)} = \frac{x^{(\ell-1)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$x^{(\ell)} = \gamma^{(\ell)} \hat{x}^{(\ell-1)} + \beta^{(\ell)}$$

【 図 1 J 】

バッチ正規化-畳み込みニューラルネットワークを用いた推論

$$\hat{x}^{(\ell-1)} = \frac{x^{(\ell-1)} - \mu_D}{\sqrt{\sigma_D^2 + \epsilon}}$$

$$x_i^{(\ell)} = \gamma^{(\ell)} \hat{x}_i^{(\ell-1)} + \beta^{(\ell)}$$

【 図 1 K 】

畳み込みニューラルネットワークを用いたバッチ正規化バックワードパス

$$\nabla_{\gamma^{(\ell)}} \mathcal{L} = \sum_{i=1}^n (\nabla_{x^{(\ell+1)}} \mathcal{L})_i \cdot \hat{x}_i^{(\ell)}$$

$$\nabla_{\beta^{(\ell)}} \mathcal{L} = \sum_{i=1}^n (\nabla_{x^{(\ell+1)}} \mathcal{L})_i$$

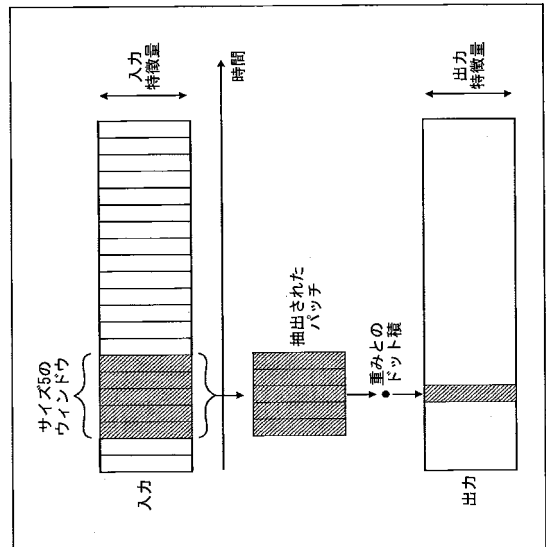
【 図 1 L 】

畳み込み層におけるバッチ正規化

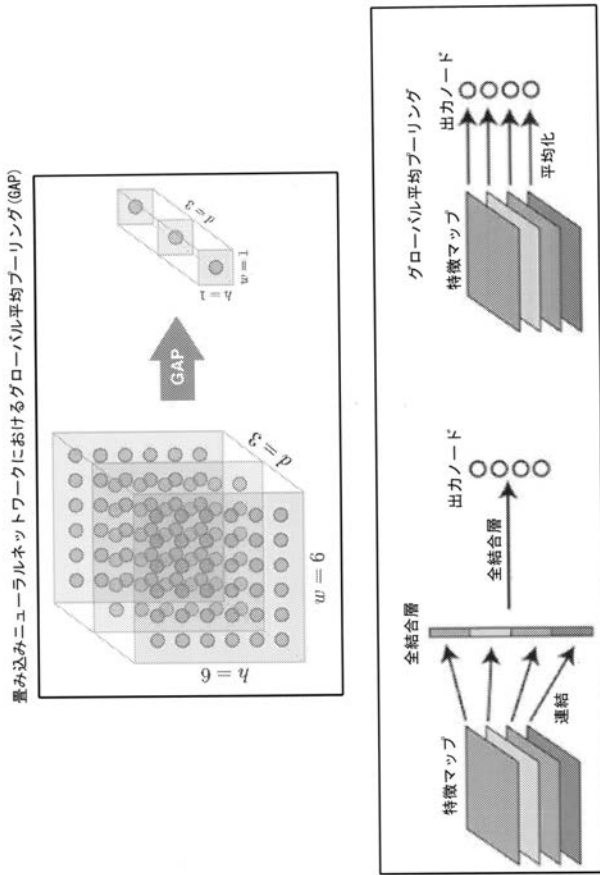
```
conv_model.add(layers.Conv2D(32, 3, activation='relu')) ← 畳み込み層の後
conv_model.add(layers.BatchNormalization())
dense_model.add(layers.Dense(32, activation='relu')) ← 密層の後
dense_model.add(layers.BatchNormalization())
```

【 図 1 M 】

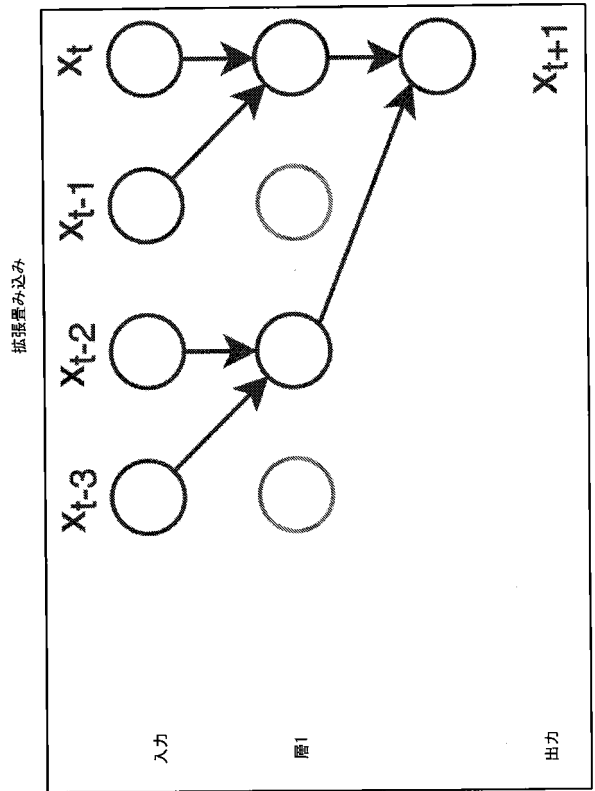
畳み込みニューラルネットワークにおいて使用されるID畳み込み



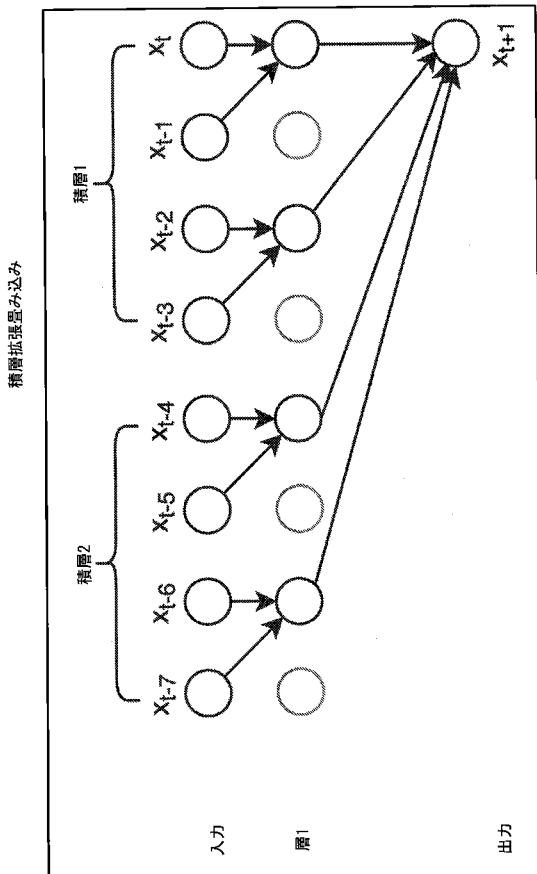
【図 1 N】



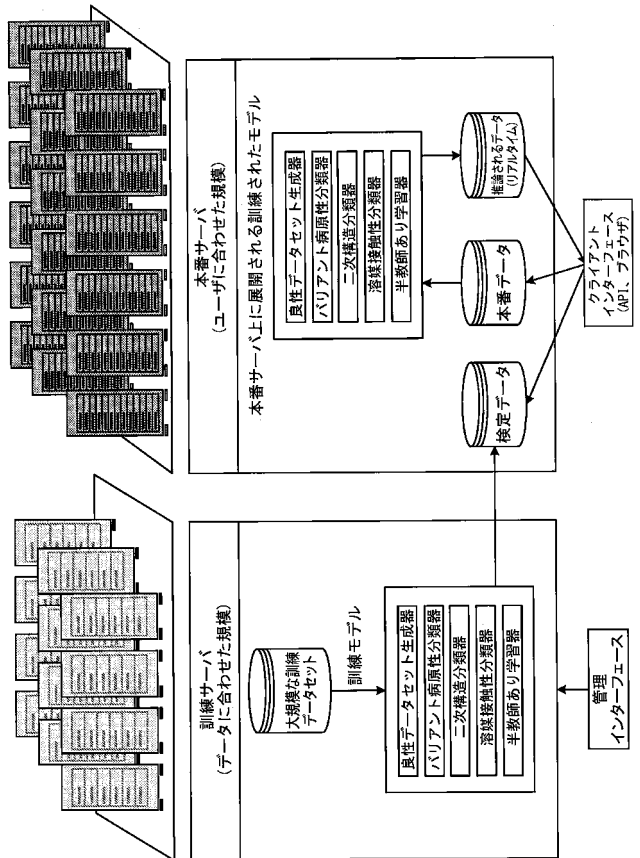
【図 1 O】



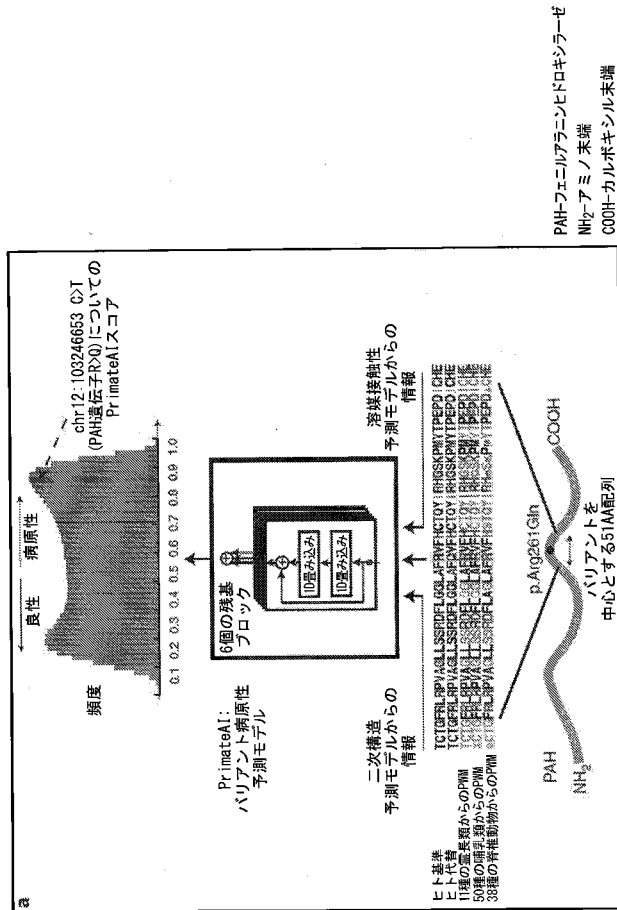
【図 1 P】



【図 1 Q】



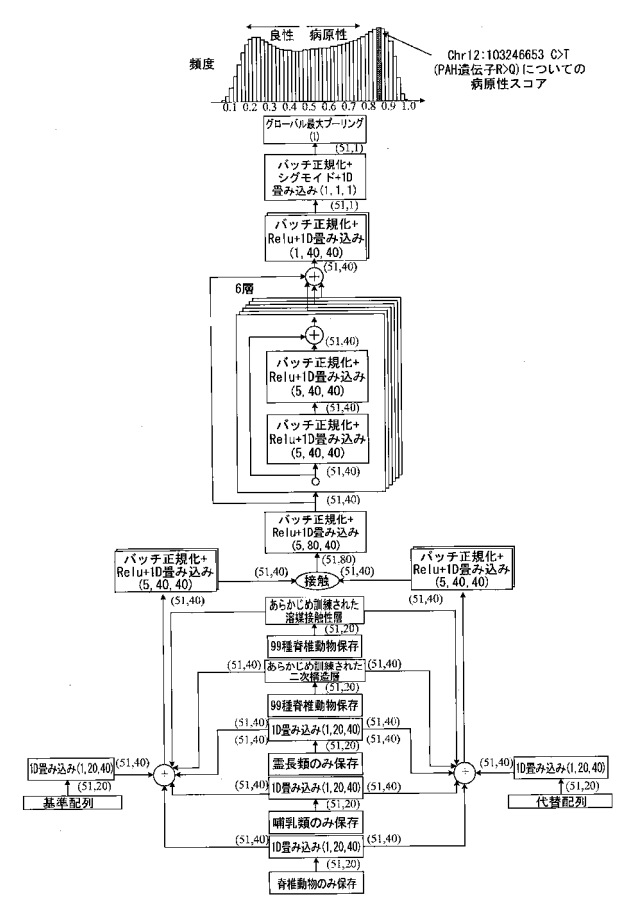
【 図 2 】



【 図 4 A 】

層	層への入力	タイプ	カーネルの数 ウィンドウサイズ	形状	活性化
層1a	基準配列	畳み込み1D	40,1	(51,40)	線形
層1b	代替配列	畳み込み1D	40,1	(51,40)	線形
層1c	畜長類保存	畳み込み1D	40,1	(51,40)	線形
層1d	哺乳類保存	畳み込み1D	40,1	(51,40)	線形
層1e	脊椎動物保存	畳み込み1D	40,1	(51,40)	線形
層1f	あらかじめ訓練された 二次構造予測モデル (訓練可能ネットワーク)	畳み込み1D	40,1	(51,40)	線形
層1g	あらかじめ訓練された 溶媒接触性モデル (訓練可能ネットワーク)	畳み込み1D	40,1	(51,40)	線形
マージ1a	層1a、層1c、層1d、 層1e、層1f、層1g	加算	-	(51,40)	-
マージ1b	層1b、層1c、層1d、 層1e、層1f、層1g	加算	-	(51,40)	-
層2a	マージ1a	バッチ正規化	-	(51,40)	-
層3a	層2a	活性化	-	(51,40)	線形
層4a	層3a	畳み込み1D	40,5	(51,40)	線形
層5a	層4a	バッチ正規化	-	(51,40)	-
層6a	層5a	活性化	-	(51,40)	線形
層7a	層6a	畳み込み1D	40,5	(51,40)	線形
層2b	マージ1b	バッチ正規化	-	(51,40)	-
層3b	層2b	活性化	-	(51,40)	線形
層4b	層3b	畳み込み1D	40,5	(51,40)	線形

【 図 3 】



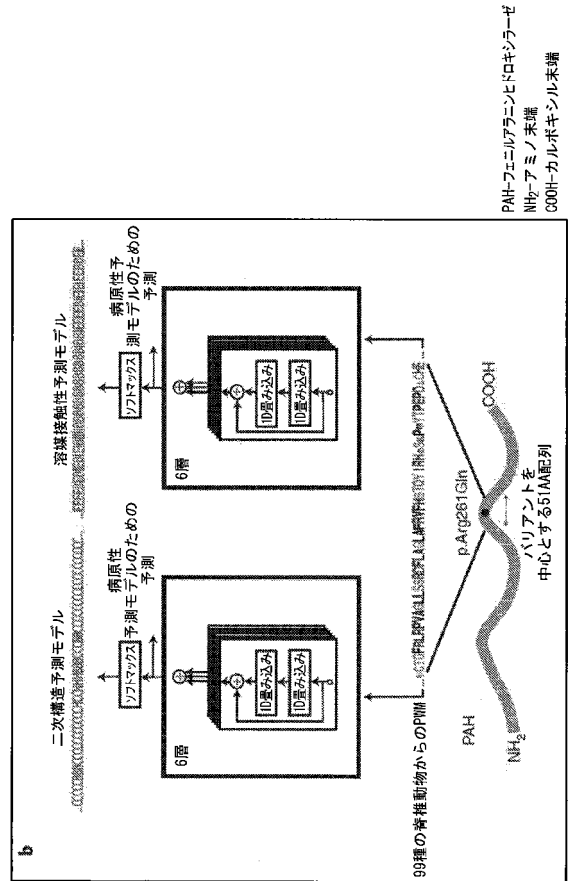
【 図 4 B 】

層5b	層4b	バッチ正規化	-	(51,40)	-
層6b	層5b	活性化	-	(51,40)	ReLU
層7b	層6b	畳み込み1D	40,5	(51,40)	線形
マージ2a	層7a、層7b	連結	-	(51,80)	-
マージ2b	層7a、層7b	連結	-	(51,80)	-
層8a	マージ2a	畳み込み1D	40,5	(51,40)	線形
層8b	マージ2b	畳み込み1D	40,5	(51,40)	線形
層9	層8a	バッチ正規化	-	(51,40)	-
層10	層9	活性化	-	(51,40)	ReLU
層11	層10	畳み込み1D	40,5	(51,40)	線形
層12	層11	バッチ正規化	-	(51,40)	-
層13	層12	活性化	-	(51,40)	ReLU
層14	層13	畳み込み1D	40,5	(51,40)	線形
マージ3	層8a、層14	加算	-	(51,40)	-
層15	マージ3	バッチ正規化	-	(51,40)	-
層16	層15	活性化	-	(51,40)	ReLU
層17	層16	畳み込み1D	40,5	(51,40)	線形
層18	層17	バッチ正規化	-	(51,40)	-
層19	層18	活性化	-	(51,40)	ReLU
層20	層19	畳み込み1D	40,5	(51,40)	線形
マージ4	マージ3、層20	加算	-	(51,40)	-
層21	マージ4	畳み込み1D	40,1	(51,40)	線形
層22	マージ4	バッチ正規化	-	(51,40)	-
層23	層22	活性化	-	(51,40)	ReLU
層24	層23	畳み込み1D	40,5	(51,40)	線形
層25	層24	バッチ正規化	-	(51,40)	-
層26	層25	活性化	-	(51,40)	ReLU
層27	層26	畳み込み1D	40,5	(51,40)	線形
マージ5	マージ4、層27	加算	-	(51,40)	-
層28	マージ5	バッチ正規化	-	(51,40)	-
層29	層28	活性化	-	(51,40)	ReLU
層30	層29	畳み込み1D	40,5	(51,40)	線形
層31	層30	バッチ正規化	-	(51,40)	-
層32	層31	活性化	-	(51,40)	ReLU

【図4C】

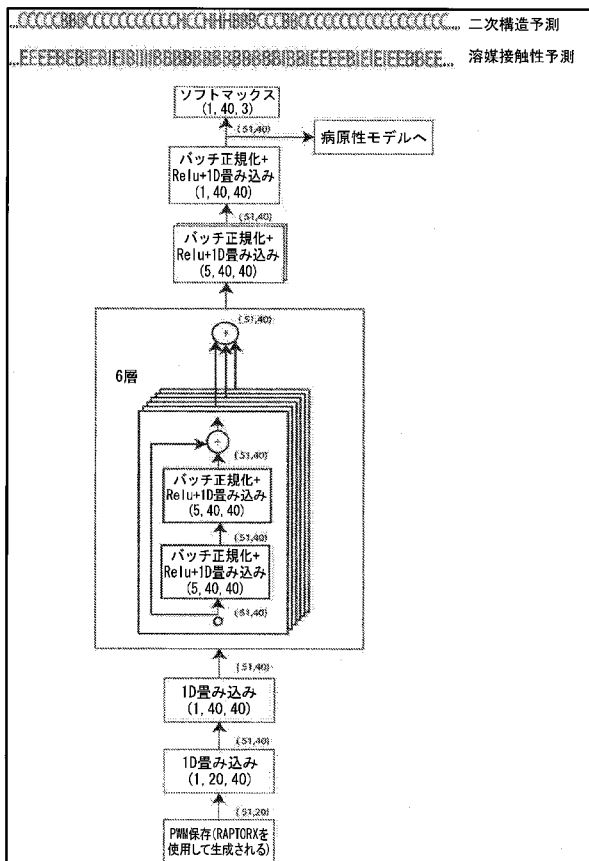
層33	層32	畳み込みID	40.5	(51,40)	線形
マージ6	マージ5、層33	加算	-	(51,40)	-
層34	マージ6	畳み込みID	40.1	(51,40)	線形
層35	マージ6	バッチ正規化	-	(51,40)	-
層36	層35	活性化	-	(51,40)	Relu
層37	層36	畳み込みID	40.5	(51,40)	線形
層38	層37	バッチ正規化	-	(51,40)	-
層39	層38	活性化	-	(51,40)	Relu
層40	層39	畳み込みID	40.5	(51,40)	線形
マージ7	マージ6、層40	加算	-	(51,40)	-
層41	マージ7	バッチ正規化	-	(51,40)	-
層42	層41	活性化	-	(51,40)	Relu
層43	層42	畳み込みID	40.5	(51,40)	線形
層44	層43	バッチ正規化	-	(51,40)	-
層45	層44	活性化	-	(51,40)	Relu
層46	層45	畳み込みID	40.5	(51,40)	線形
マージ8	マージ7、層46	加算	-	(51,40)	-
層47	マージ8	畳み込みID	40.1	(51,40)	線形
マージ9	層4b、層21、層34、層47	加算	-	(51,40)	-
層48	マージ10	バッチ正規化	-	(51,40)	-
層49	層48	活性化	-	(51,40)	Relu
層50	層49	畳み込みID	40.1	(51,40)	線形
層51	層50	バッチ正規化	-	(51,40)	-
層52	層51	活性化	-	(51,40)	Relu
層53	層52	畳み込みID	40.1	(51,40)	線形
マージ10	マージ9、層53	加算	-	(51,40)	-
層54	マージ10	畳み込みID	1.1	(51.1)	シグモイド
出力層	層54	グローバル最大プーリングID	-	-	-

【図5】



PAH-フェニルアラニンヒドロキシラーゼ  
NH<sub>2</sub>-アミノ末端  
COOH-カルボキシル基末端

【図6】



【図7A】

層	層への入力	タイプ	カーネルの数&カーネルウィンドウサイズ	形状	活性化
層1a	入力PSFM	畳み込みID	40.1	(51,40)	線形
層2a	層1a	畳み込みID	40.1	(51,40)	線形
層1b	入力PSFM	畳み込みID	40.1	(51,40)	線形
層2b	層1b	畳み込みID	40.1	(51,40)	線形
層3	層2a	バッチ正規化	-	(51,40)	-
層4	層3	活性化	-	(51,40)	Relu
層5	層4	畳み込みID	40.5	(51,40)	線形
層6	層5	バッチ正規化	-	(51,40)	-
層8	層7	活性化	-	(51,40)	Relu
層9	層8	畳み込みID	40.5	(51,40)	線形
マージ1	層2a、層9	加算	-	(51,40)	-
層10	マージ1	バッチ正規化	-	(51,40)	-
層11	層10	活性化	-	(51,40)	Relu
層12	層11	畳み込みID	40.5	(51,40)	線形
層13	層12	バッチ正規化	-	(51,40)	-
層14	層13	活性化	-	(51,40)	Relu
層15	層14	畳み込みID	40.5	(51,40)	線形
マージ2	マージ1、層15	加算	-	(51,40)	-
層16	マージ2	畳み込みID	40.1	(51,40)	線形
層17	マージ2	バッチ正規化	-	(51,40)	-
層18	層17	活性化	-	(51,40)	Relu
層19	層18	畳み込みID	40.5	(51,40)	線形
層20	層19	バッチ正規化	-	(51,40)	-
層21	層20	活性化	-	(51,40)	Relu
層22	層21	畳み込みID	40.5	(51,40)	線形
マージ3	マージ2、層22	加算	-	(51,40)	-

【図7B】

層23	マージ3	バッチ正規化	-	(51.40)	-
層24	層23	活性化	-	(51.40)	Relu
層25	層24	畳み込み1D	40.5	(51.40)	線形
層26	層25	バッチ正規化	-	(51.40)	-
層27	層26	活性化	-	(51.40)	Relu
層28	層27	畳み込み1D	40.5	(51.40)	線形
マージ4	マージ3、層28	加算	-	(51.40)	-
層29	マージ4	畳み込み1D	40.1	(51.40)	線形
層30	マージ4	バッチ正規化	-	(51.40)	-
層31	層30	活性化	-	(51.40)	Relu
層32	層31	畳み込み1D	40.5	(51.40)	線形
層33	層32	バッチ正規化	-	(51.40)	-
層34	層33	活性化	-	(51.40)	Relu
層35	層34	畳み込み1D	40.5	(51.40)	線形
マージ5	マージ4、層35	加算	-	(51.40)	-
層36	マージ5	バッチ正規化	-	(51.40)	-
層37	層36	活性化	-	(51.40)	Relu
層38	層37	畳み込み1D	40.5	(51.40)	線形
層39	層38	バッチ正規化	-	(51.40)	-
層40	層39	活性化	-	(51.40)	Relu
層41	層40	畳み込み1D	40.5	(51.40)	線形
マージ6	マージ5、層41	加算	-	(51.40)	-
層42	マージ6	畳み込み1D	40.1	(51.40)	線形
マージ7	層2b、層16、層29、層42	加算	-	(51.40)	-
層43	マージ7	バッチ正規化	-	(51.40)	-
層44	層43	活性化	-	(51.40)	Relu
層45	マージ44	畳み込み1D	40.1	(51.40)	線形
層46	層45	バッチ正規化	-	(51.40)	-
層47	層46	活性化	-	(51.40)	Relu
層48	層47	畳み込み1D	40.1	(51.40)	線形
マージ8	マージ7、層48	加算	-	(51.40)	-
出力層	マージ8	畳み込み1D	3.1	(51.3)	ソフトマックス

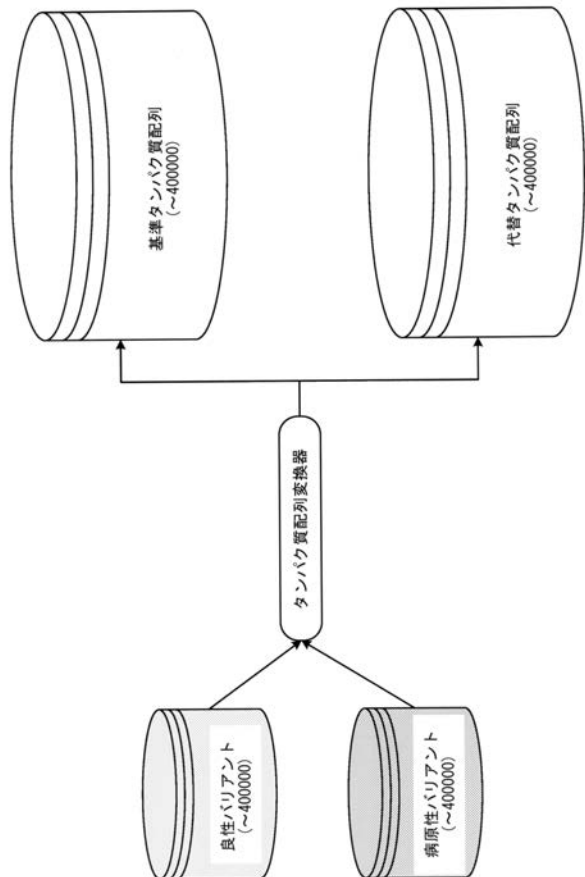
【図8A】

層	層への入力	タイプ	カーネルの数&カーネルウィンドウサイズ	形状	活性化
層1a	入力PSFM	畳み込み1D	40.1	(51.40)	線形
層1b	層1a	畳み込み1D	40.1	(51.40)	線形
層2a	層1b	畳み込み1D	40.1	(51.40)	線形
層2b	層2a	バッチ正規化	-	(51.40)	-
層3	層2b	活性化	-	(51.40)	Relu
層4	層3	畳み込み1D	40.5	(51.40)	線形
層5	層4	バッチ正規化	-	(51.40)	-
層6	層5	活性化	-	(51.40)	Relu
層7	層6	畳み込み1D	40.5	(51.40)	線形
層8	層7	バッチ正規化	-	(51.40)	-
層9	層8	活性化	-	(51.40)	Relu
マージ1	層2a、層9	加算	-	(51.40)	-
層10	マージ1	畳み込み1D	40.1	(51.40)	線形
層11	層10	バッチ正規化	-	(51.40)	-
層12	層11	活性化	-	(51.40)	Relu
層13	層12	畳み込み1D	40.5	(51.40)	線形
層14	層13	バッチ正規化	-	(51.40)	-
層15	層14	活性化	-	(51.40)	Relu
層16	層15	畳み込み1D	40.5	(51.40)	線形
マージ2	マージ1、層16	加算	-	(51.40)	-
層17	マージ2	畳み込み1D	40.1	(51.40)	線形
層18	層17	バッチ正規化	-	(51.40)	-
層19	層18	活性化	-	(51.40)	Relu
層20	層19	畳み込み1D	40.5	(51.40)	線形
層21	層20	バッチ正規化	-	(51.40)	-
層22	層21	活性化	-	(51.40)	Relu
マージ3	層22	畳み込み1D	40.5	(51.40)	線形
マージ3	マージ2、層22	加算	-	(51.40)	-

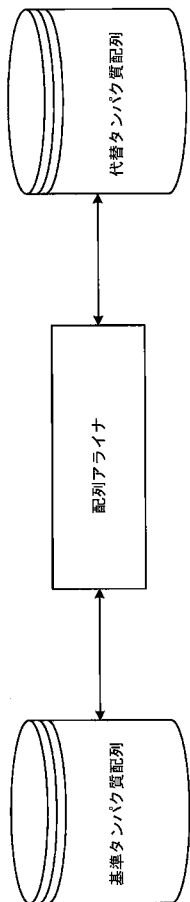
【図8B】

層23	マージ3	バッチ正規化	-	(51.40)	-
層24	層23	活性化	-	(51.40)	Relu
層25	層24	畳み込み1D	40.5	(51.40)	線形
層26	層25	バッチ正規化	-	(51.40)	-
層27	層26	活性化	-	(51.40)	Relu
層28	層27	畳み込み1D	40.5	(51.40)	線形
マージ4	マージ3、層28	加算	-	(51.40)	-
層29	マージ4	畳み込み1D	40.1	(51.40)	線形
層30	層29	バッチ正規化	-	(51.40)	-
層31	層30	活性化	-	(51.40)	Relu
層32	層31	畳み込み1D	40.5	(51.40)	線形
層33	層32	バッチ正規化	-	(51.40)	-
層34	層33	活性化	-	(51.40)	Relu
層35	層34	畳み込み1D	40.5	(51.40)	線形
マージ5	マージ4、層35	加算	-	(51.40)	-
層36	マージ5	バッチ正規化	-	(51.40)	-
層37	層36	活性化	-	(51.40)	Relu
層38	層37	畳み込み1D	40.5	(51.40)	線形
層39	層38	バッチ正規化	-	(51.40)	-
層40	層39	活性化	-	(51.40)	Relu
層41	層40	畳み込み1D	40.5	(51.40)	線形
マージ6	マージ5、層41	加算	-	(51.40)	-
層42	マージ6	畳み込み1D	40.1	(51.40)	線形
マージ7	層2b、層16、層29、層42	加算	-	(51.40)	-
層43	マージ7	バッチ正規化	-	(51.40)	-
層44	層43	活性化	-	(51.40)	Relu
層45	マージ44	畳み込み1D	40.1	(51.40)	線形
層46	層45	バッチ正規化	-	(51.40)	-
層47	層46	活性化	-	(51.40)	Relu
層48	層47	畳み込み1D	40.1	(51.40)	線形
マージ8	マージ7、層48	加算	-	(51.40)	-
出力層	マージ8	畳み込み1D	3.1	(51.3)	ソフトマックス

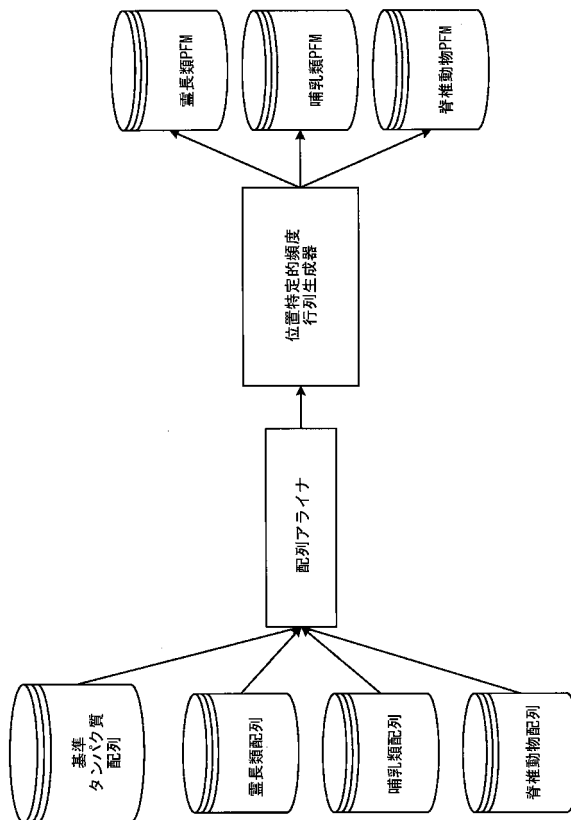
【図9】



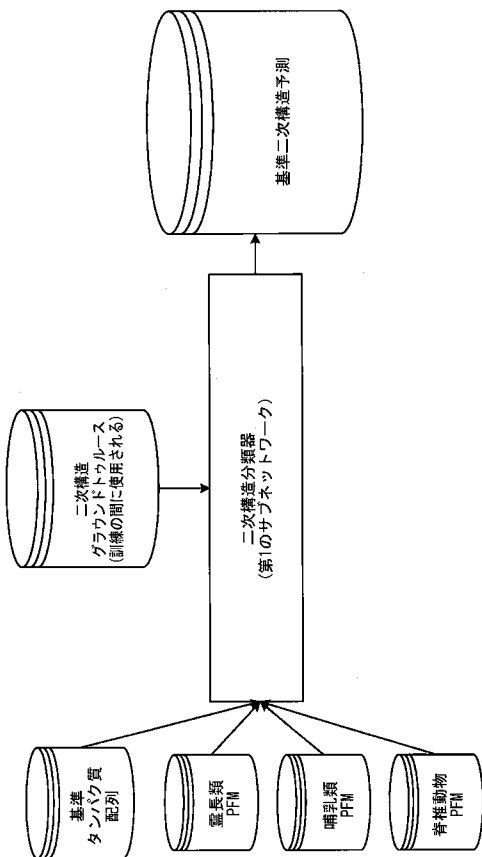
【図10】



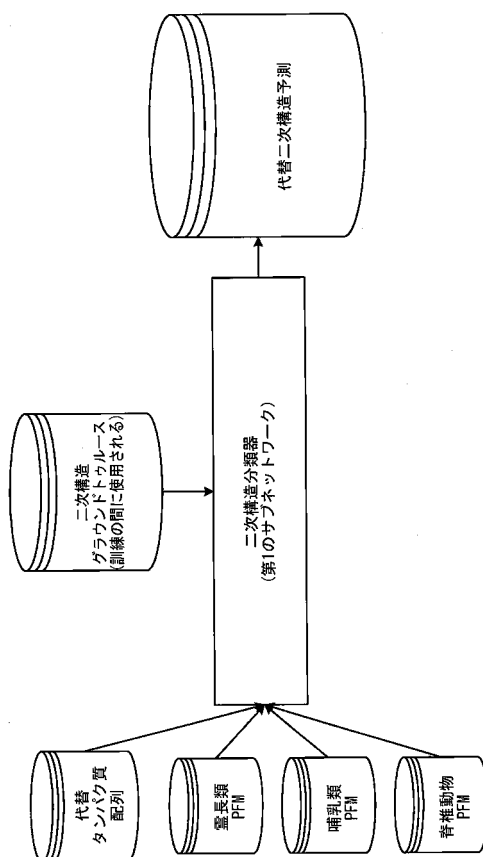
【図11】



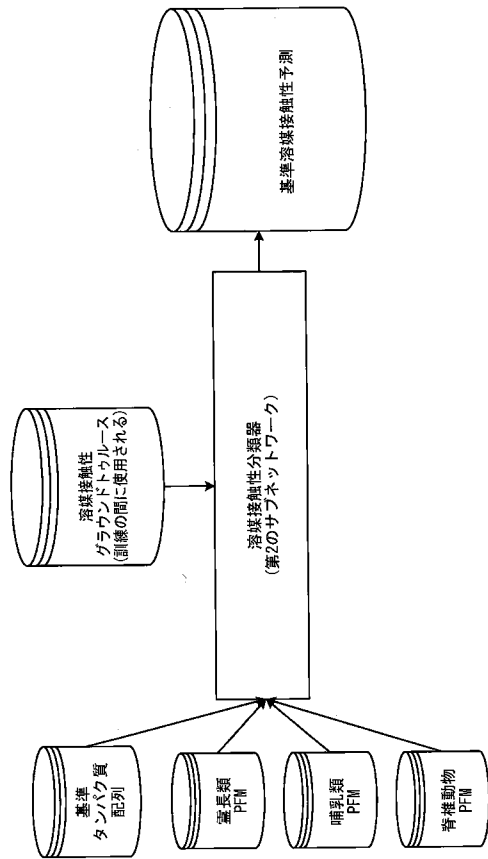
【図12】



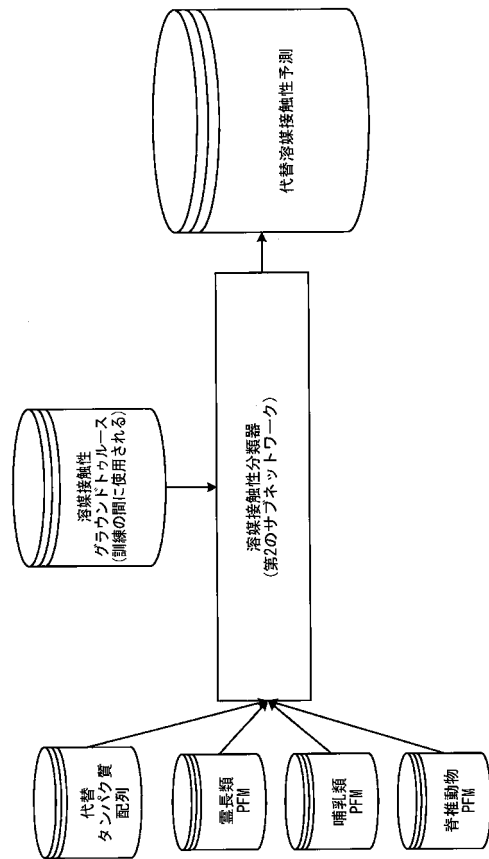
【図13】



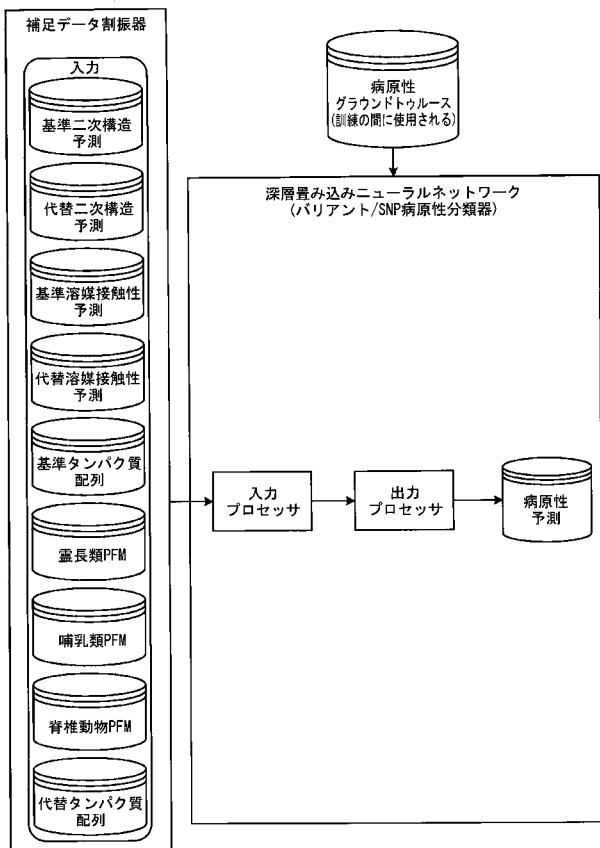
【 図 1 4 】



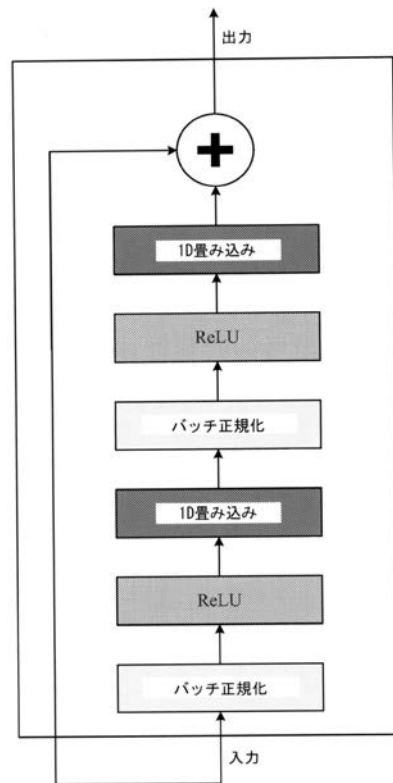
【 図 1 5 】



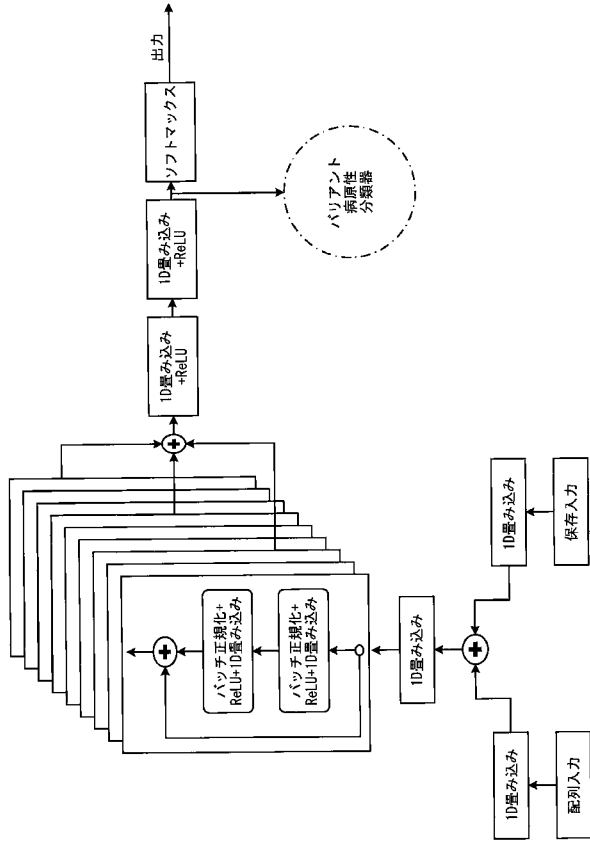
【 図 1 6 】



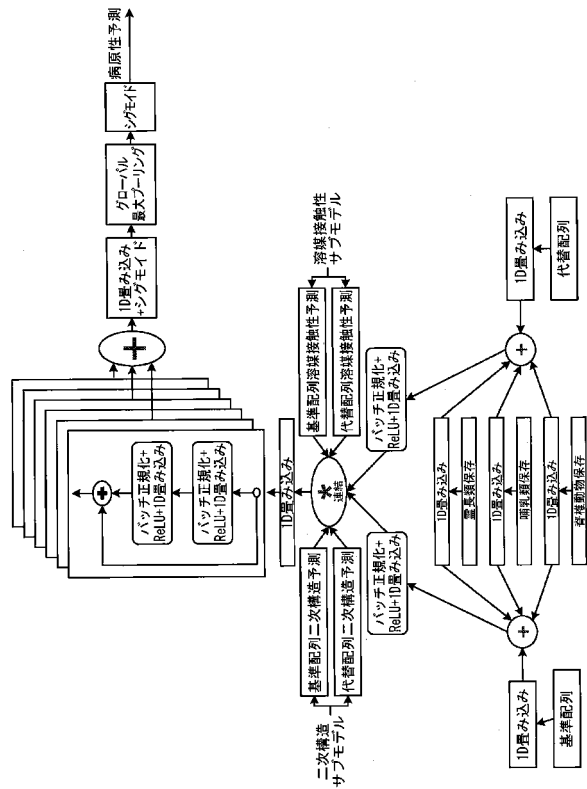
【 図 1 7 】



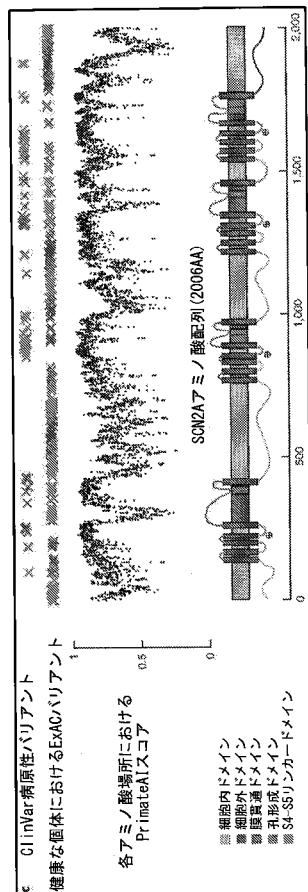
【 図 18 】



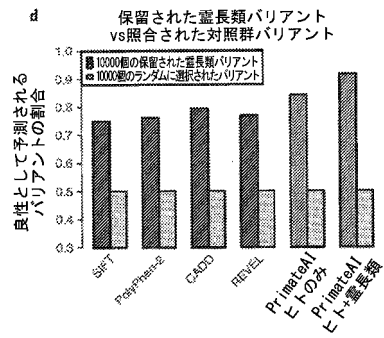
【 図 19 】



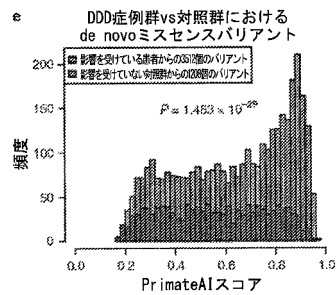
【 図 20 】



【 図 21 D 】

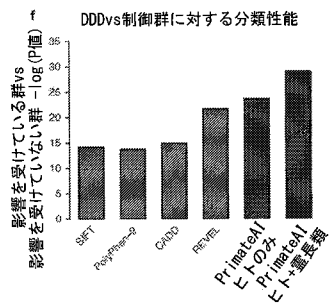


【 図 21 E 】

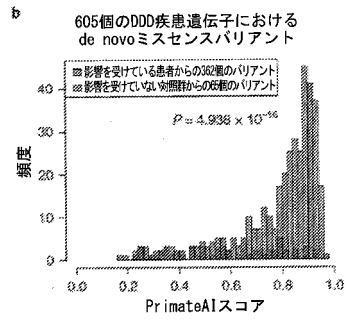




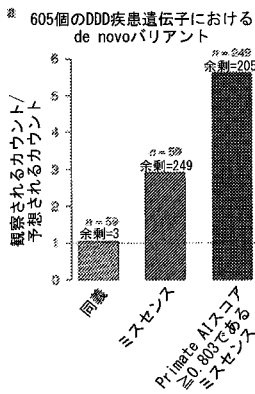
【 図 2 1 F 】



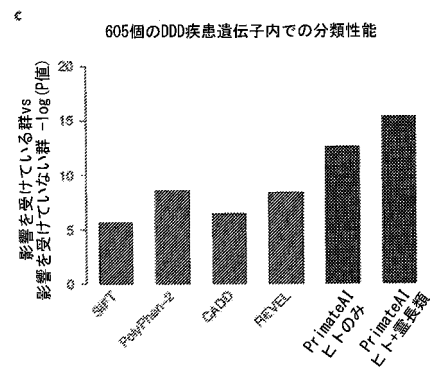
【 図 2 2 B 】



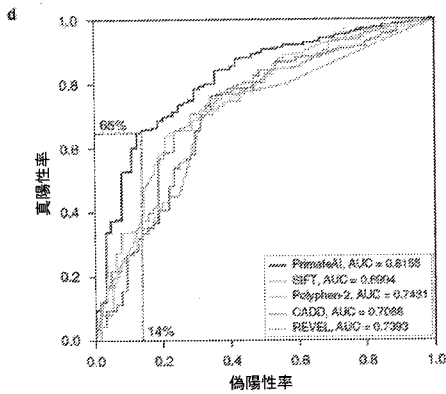
【 図 2 2 A 】



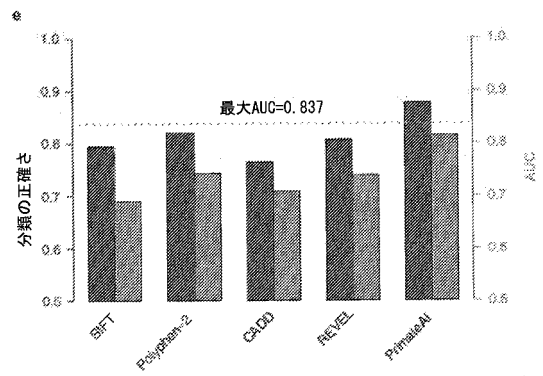
【 図 2 2 C 】



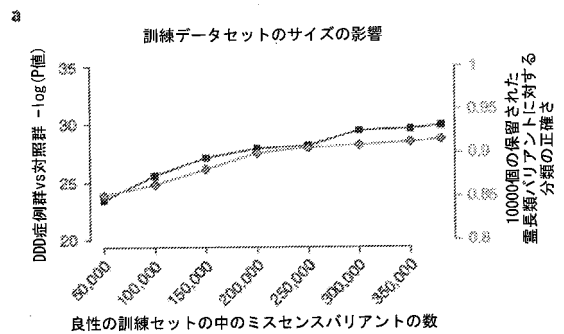
【 図 2 2 D 】



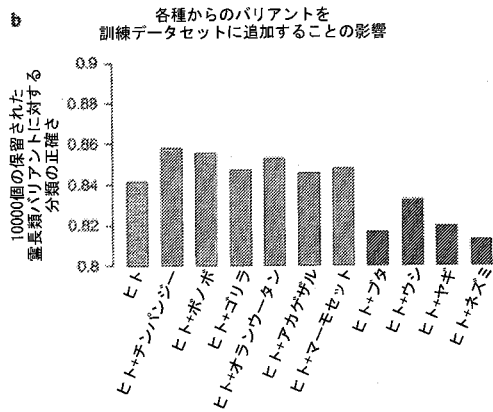
【 図 2 2 E 】



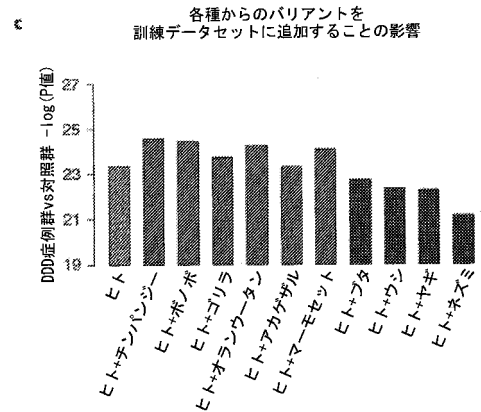
【 図 2 3 A 】



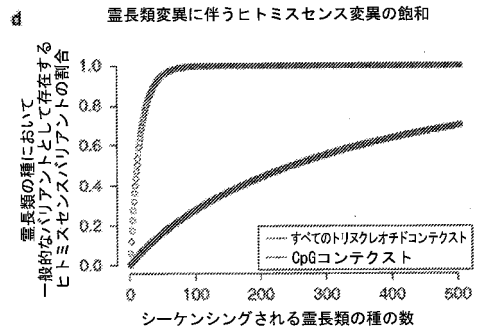
【図 2 3 B】



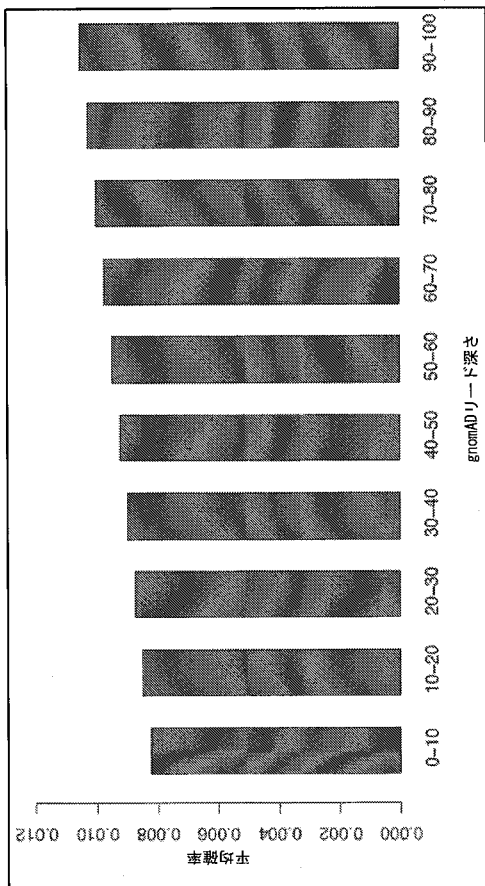
【図 2 3 C】



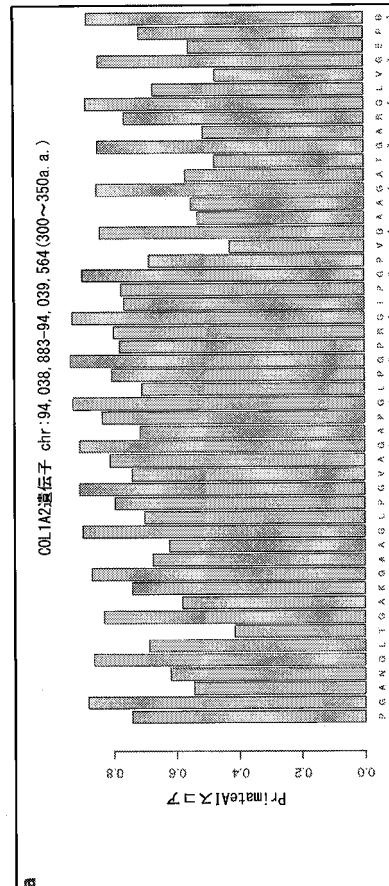
【図 2 3 D】



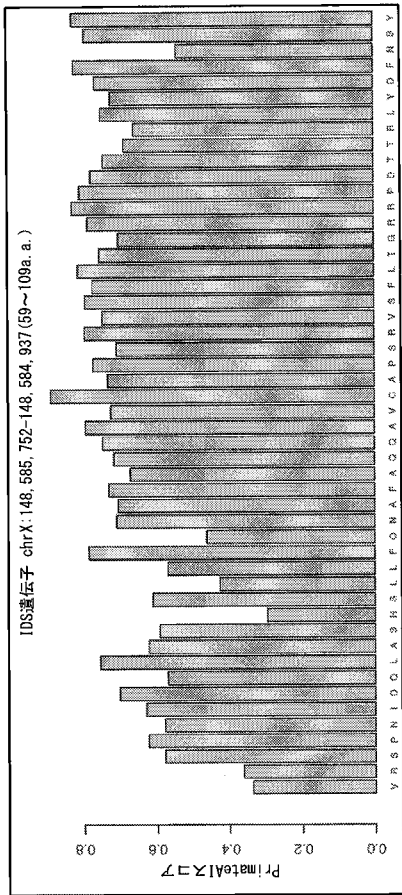
【図 2 4】



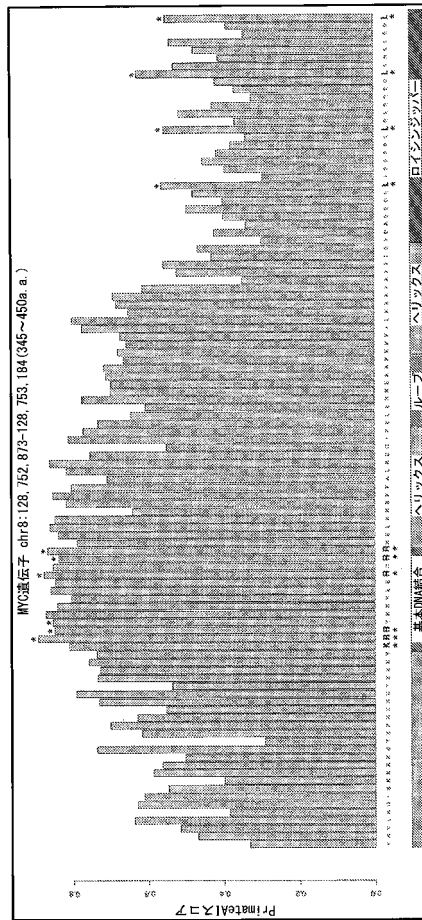
【図 2 5 A】



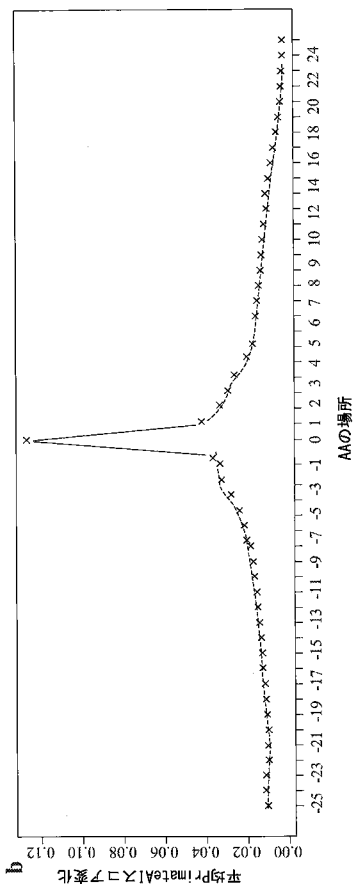
【 図 2 5 B 】



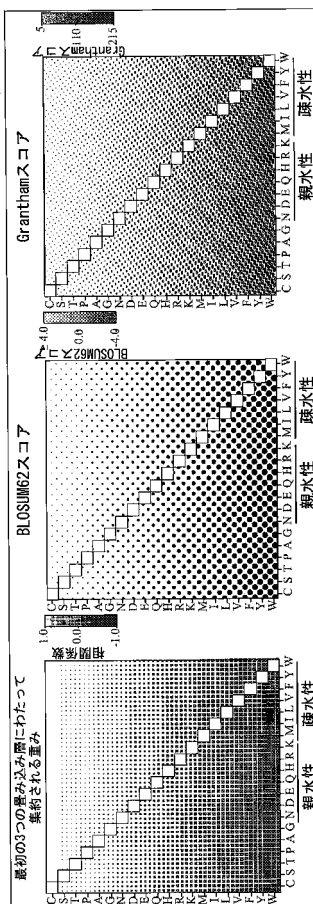
【 図 2 5 C 】



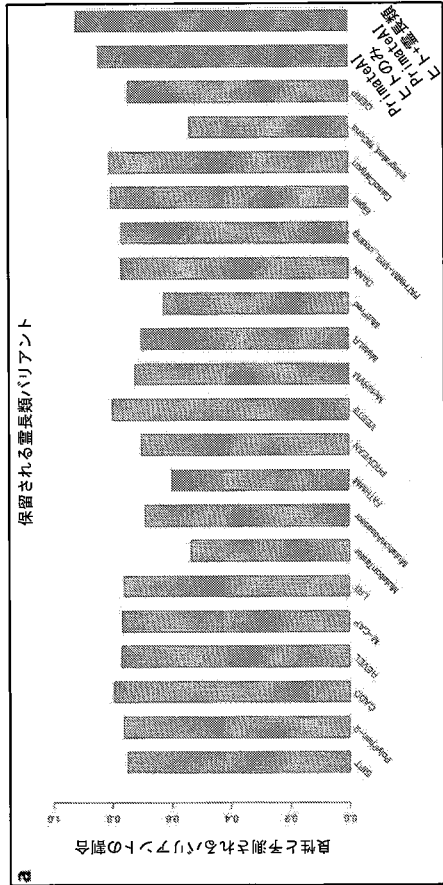
【 図 2 6 】



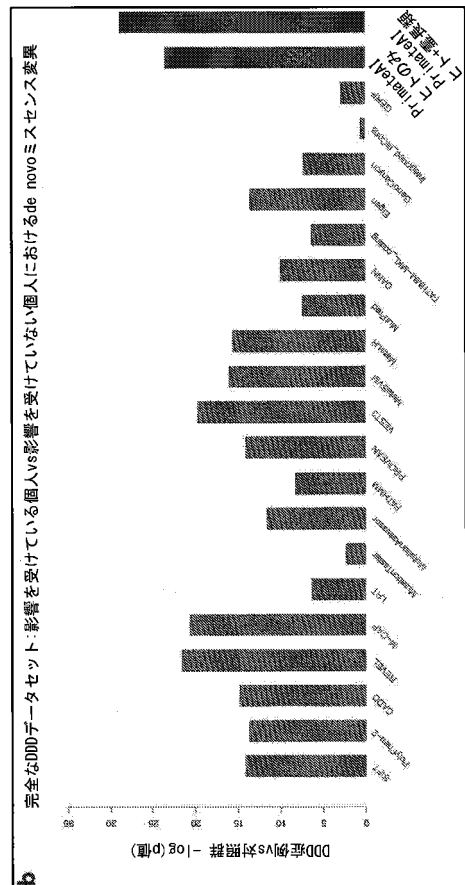
【 図 2 7 】



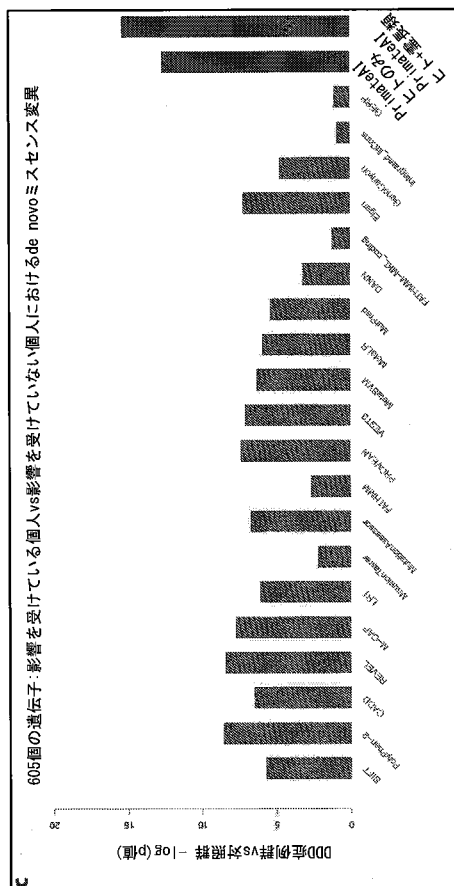
【図28A】



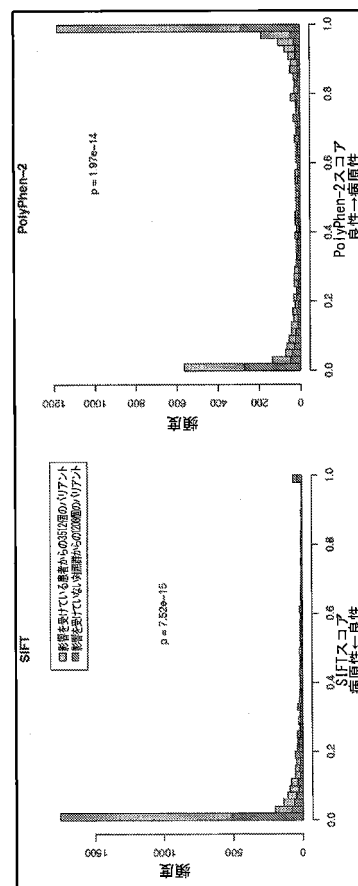
【図28B】



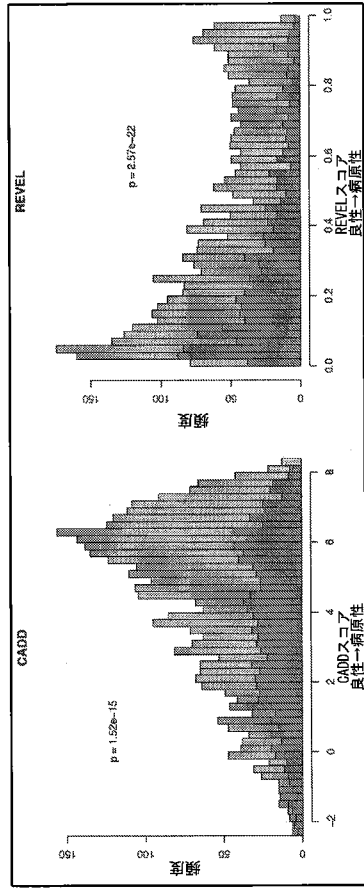
【図28C】



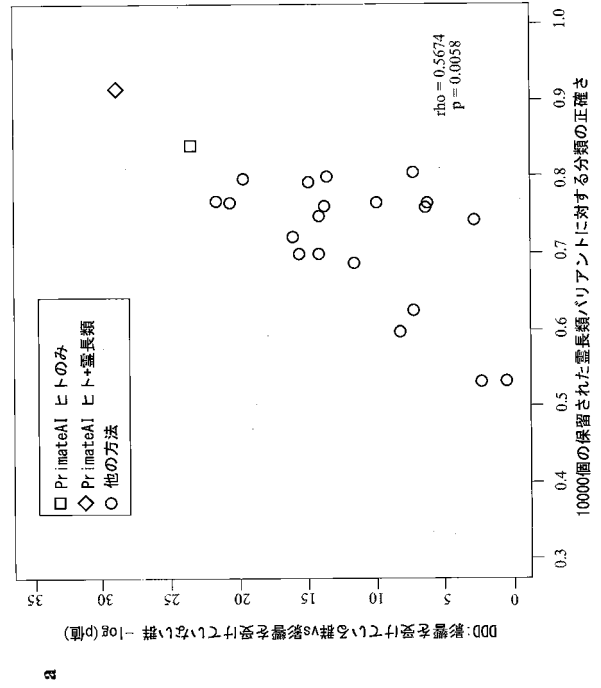
【図29A】



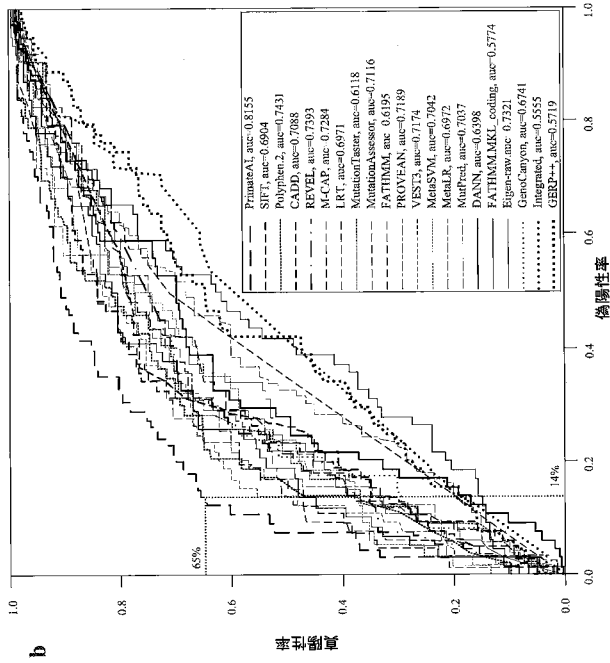
【図 29 B】



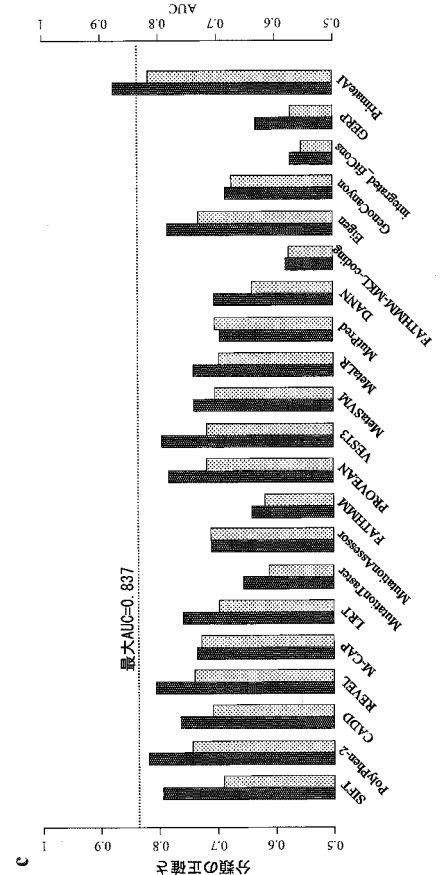
【図 30 A】



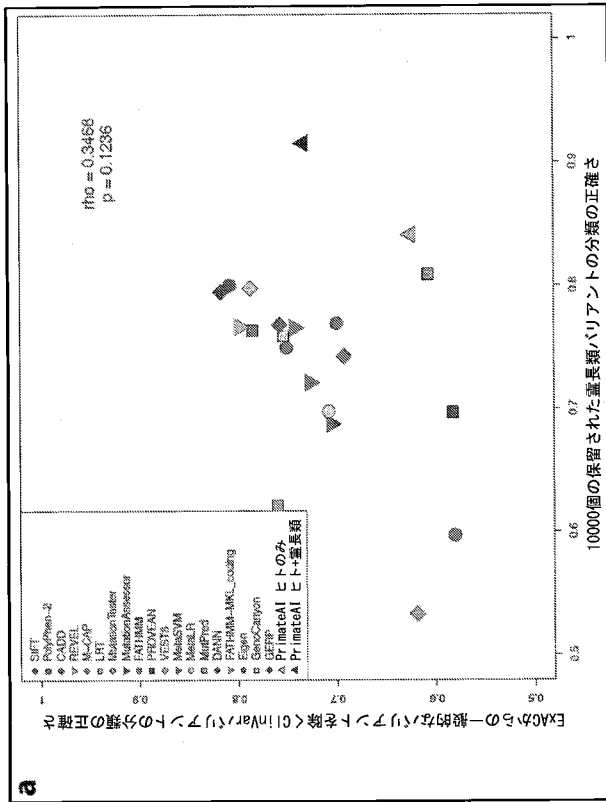
【図 30 B】



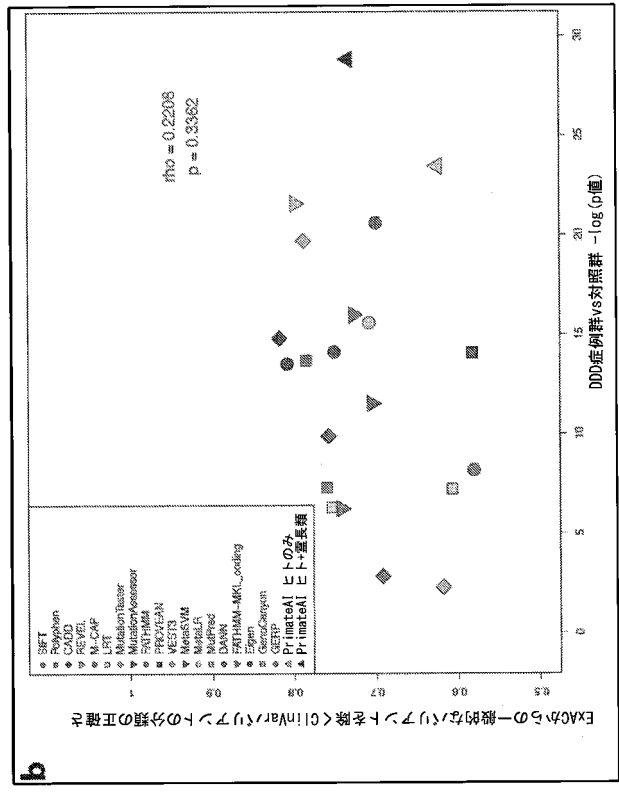
【図 30 C】



【 図 3 1 A 】



【 図 3 1 B 】



【 図 3 2 】

	訓練の正確さ	妥当性確認の正確さ	検定の正確さ
3状態二次構造	80.32 %	79.17 %	79.86 %
3状態溶解性	64.83 %	60.53 %	60.31 %

【 図 3 3 】

モデル	反復	保留された LTVの 正確さ	DDD症例群vs対照群 -log(p値)	DDD 605個の 疾患遺伝子 -log(p値)
DL予測される 二次構造レベルを 使用する	1	0.914	28.60	15.61
	2	0.919	32.32	15.94
	3	0.915	29.48	16.38
	4	0.917	30.35	15.62
	5	0.916	29.11	16.22
	中央値	0.916	29.48	15.94
可能なときに DL予測される 二次構造レベルを ヒトタンパク質の DSSP二次構造レベルで 置き換える	1	0.916	29.48	15.94
	2	0.913	31.72	16.24
	3	0.915	31.09	15.69
	4	0.915	30.49	16.09
	5	0.915	30.57	14.79
	中央値	0.915	30.79	16.09

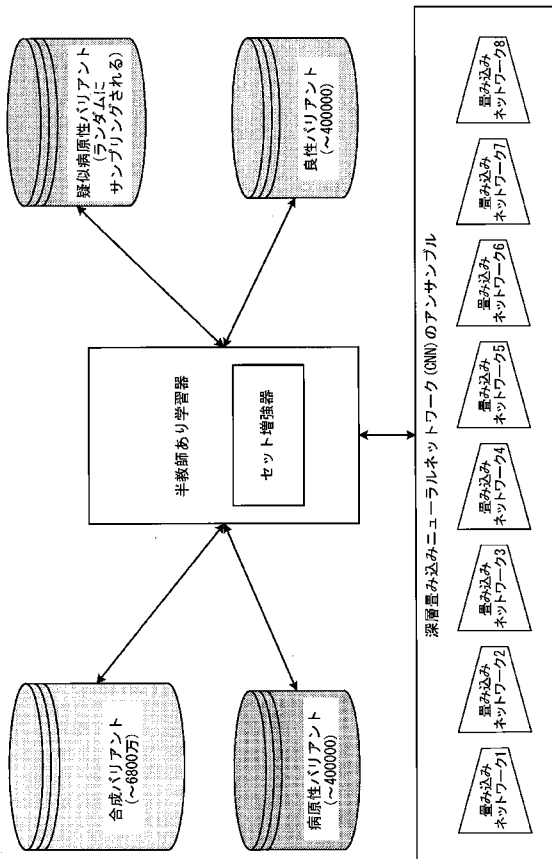
【 図 3 4 】

	保留された霊長類 バリエントの正確さ	DDD症例群vs対照群 -log(p値)
SIFT	0.7490	14.1233
PolyPhen-2	0.7618	13.7044
CADD	0.7939	14.8167
M-CAP	0.7667	20.6078
LRT	0.7598	6.2960
MutationTaster	0.5332	2.2811
MutationAssessor	0.6879	11.5202
FATHMM	0.5981	8.1690
PROVEAN	0.6995	14.0747
VEST3	0.7978	19.7068
MetaSVM	0.7215	15.9829
MetaLR	0.6991	15.5842
REVEL	0.7686	21.5892
MutPred	0.6240	7.3044
DANN	0.7666	9.8979
MKL_coding	0.7657	6.2469
Eigen	0.8003	13.5364
GenoCanyon	0.8058	7.2196
integrated_fitCons	0.5331	0.4954
GERP	0.7430	2.8285
PrimateAI ヒト	0.8411	23.4536
PrimateAI ヒト+霊長類	0.9156	28.8346

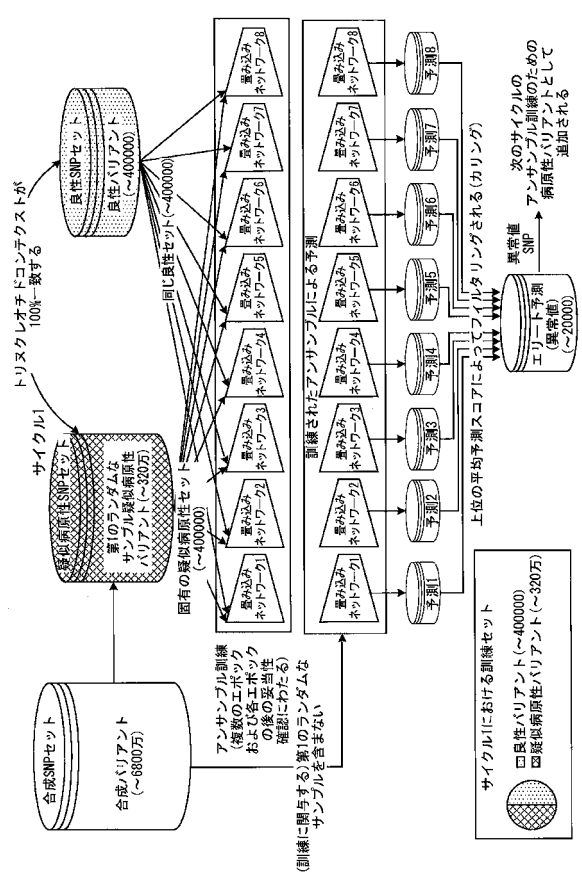
【 図 3 5 】

	-log(p値)	AUC	閾値	分類の正確さ
SIFT	5.6793	0.6903	0.001	0.7949
PolyPhen-2	8.5522	0.7430	0.956	0.8197
CADD	6.4772	0.7087	5.040	0.7643
M-CAP	7.7143	0.7283	0.081	0.7348
LRT	6.0285	0.6970	1.00E-06	0.7604
MutationTaster	2.1366	0.6117	1.000	0.6564
MutationAssessor	6.6779	0.7115	2.240	0.7096
FATHMM	2.5979	0.6194	-1.260	0.6401
PROVEAN	7.3447	0.7188	-4.310	0.7833
VEST3	7.0357	0.7173	0.679	0.7953
MetaSVM	6.2556	0.7041	-0.3371	0.7406
MetaLR	5.8707	0.6971	0.349	0.7407
REVEL	8.4115	0.7392	0.456	0.8056
MutPred	5.3370	0.7036	0.526	0.6949
DANN	3.1588	0.6397	0.996	0.7046
MKL_coding	1.1649	0.5773	0.965	0.5840
Eigen	7.1659	0.7330	0.577	0.7844
GenoCanyon	4.6763	0.6740	0.999	0.6875
integrated_fitCons	0.8316	0.5554	0.706	0.5736
GERP	1.0330	0.5718	5.060	0.6330
PrimateAI ヒト+霊長類	15.3064	0.8154	0.802	0.8772

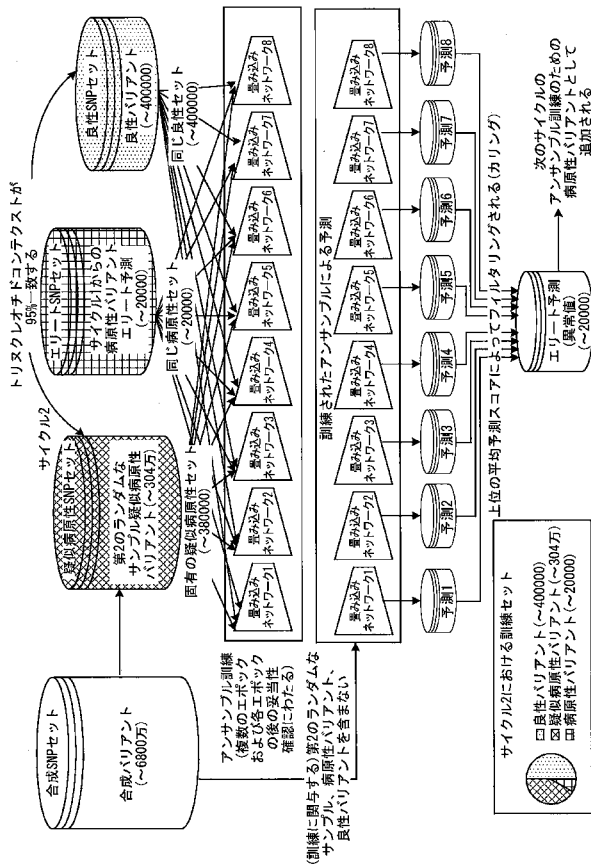
【 図 3 6 】



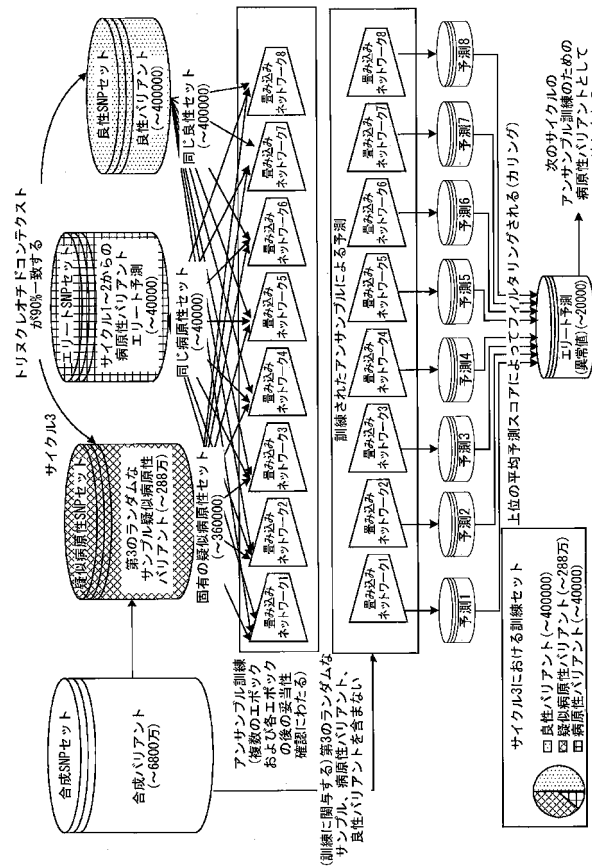
【 図 3 7 】



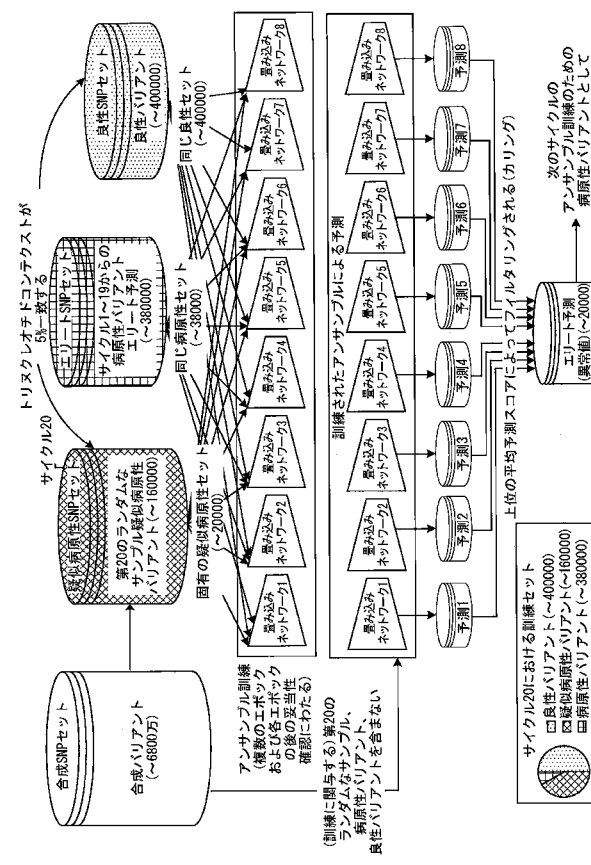
【 図 3 8 】



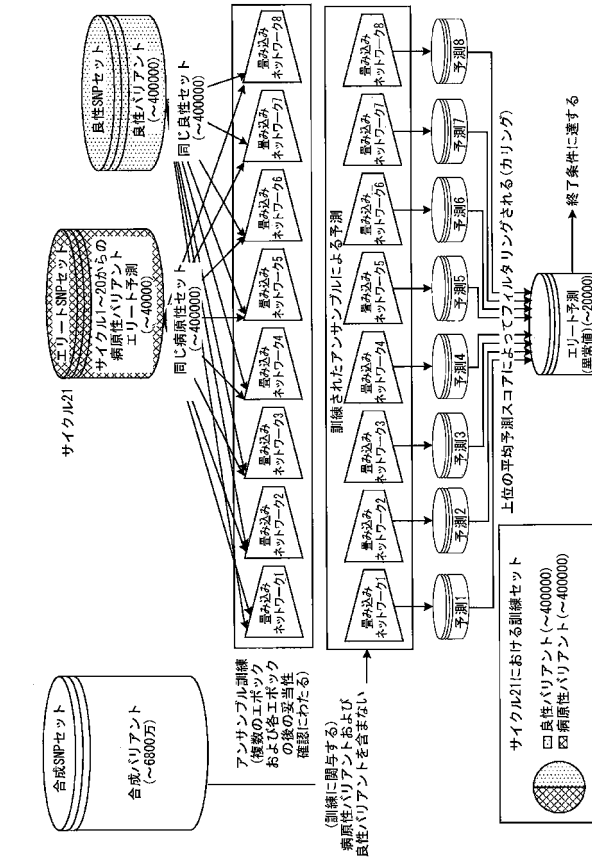
【 図 3 9 】



【 図 4 0 】

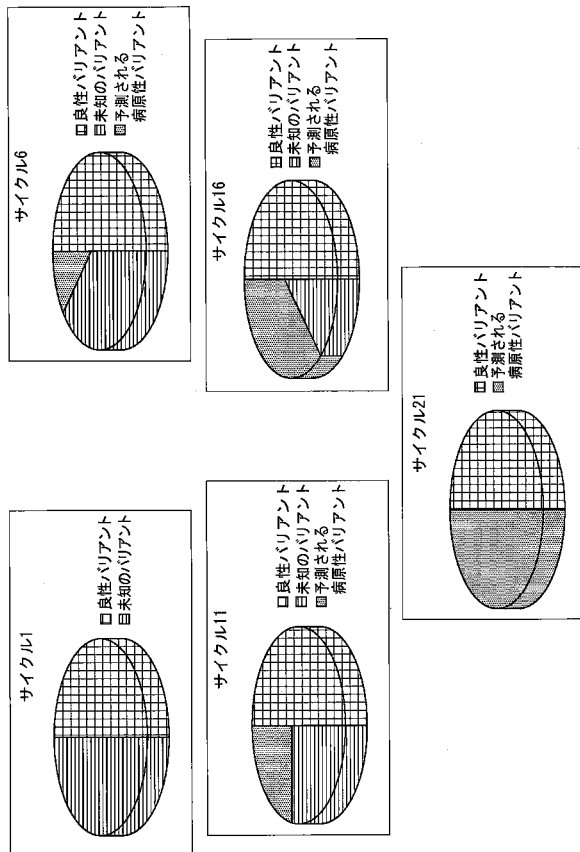


【 図 4 1 】

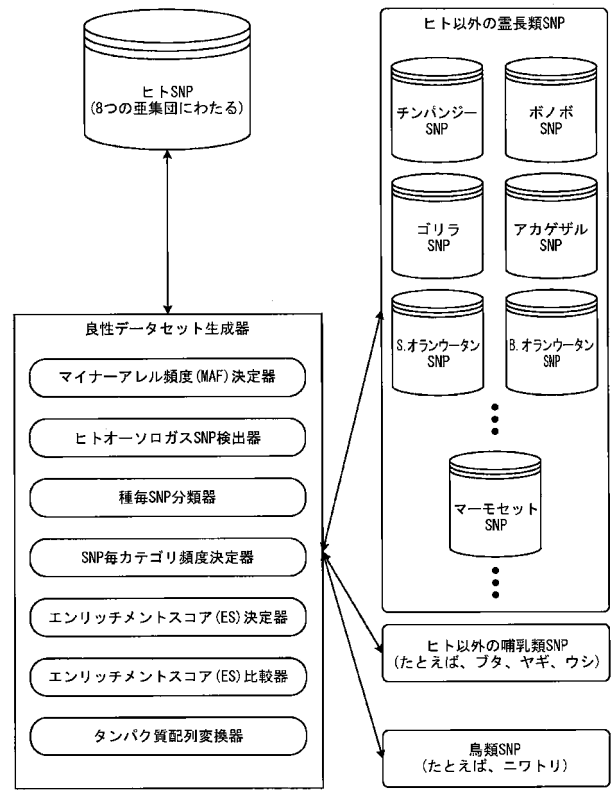




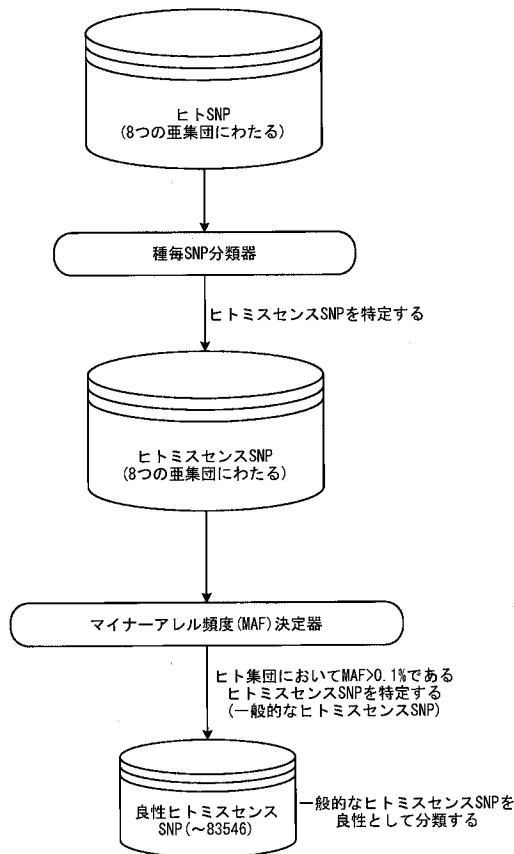
【 図 4 2 】



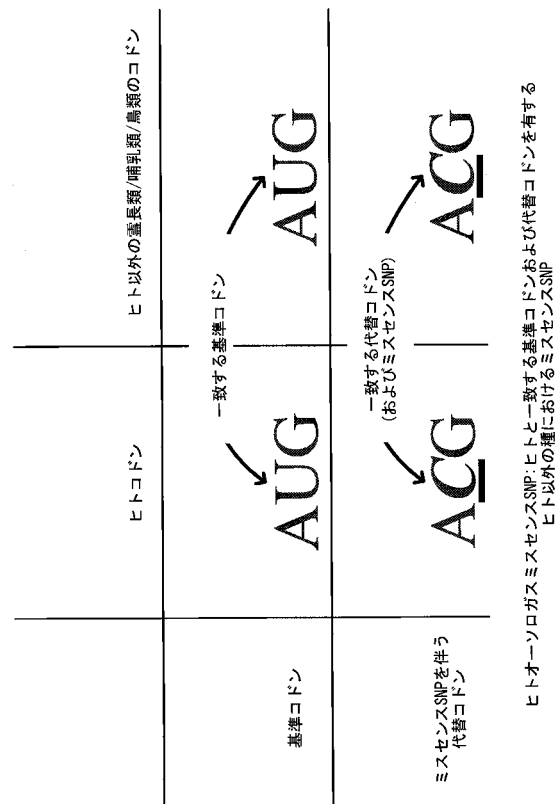
【 図 4 3 】



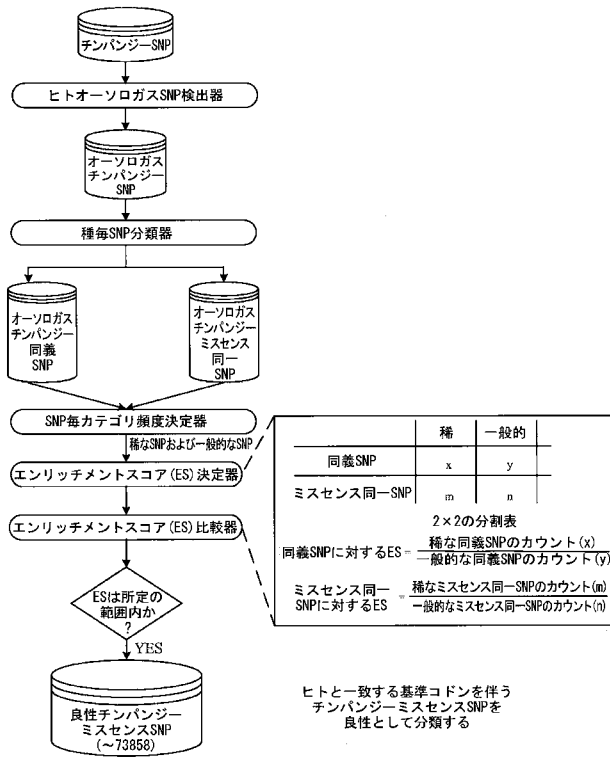
【 図 4 4 】



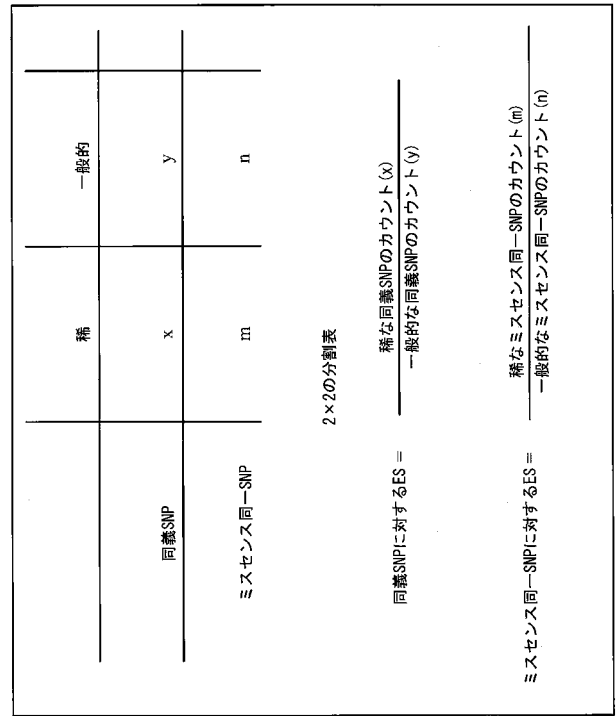
【 図 4 5 】



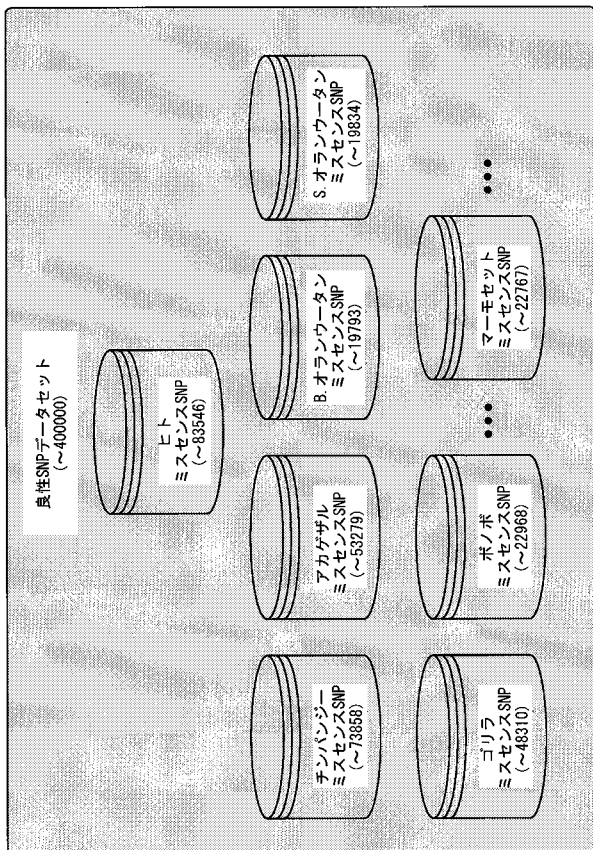
【 図 4 6 】



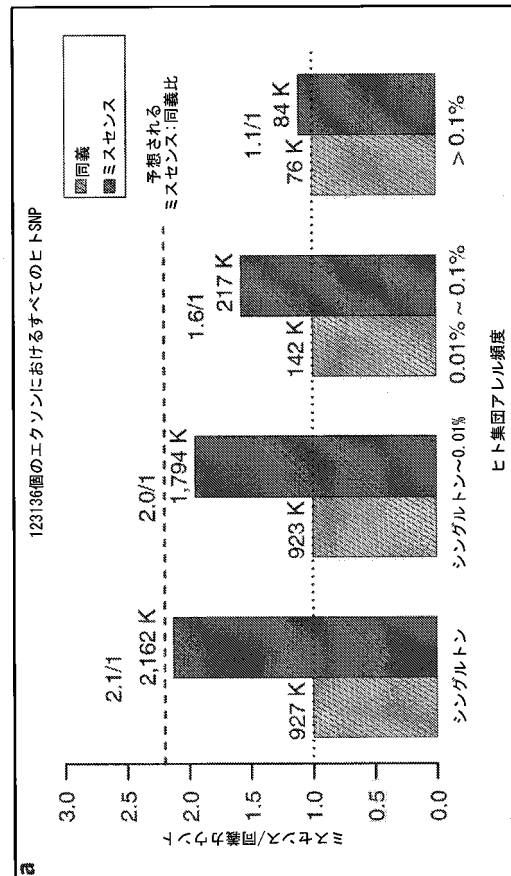
【 図 4 7 】



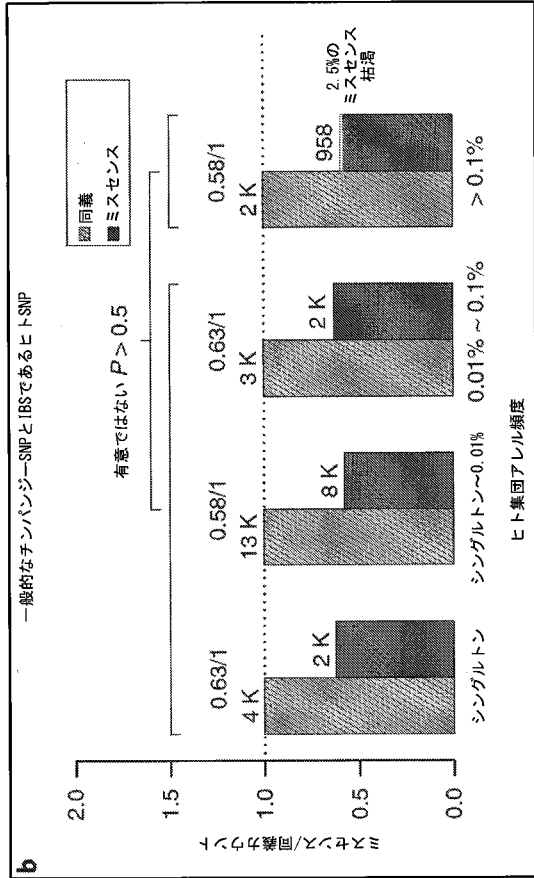
【 図 4 8 】



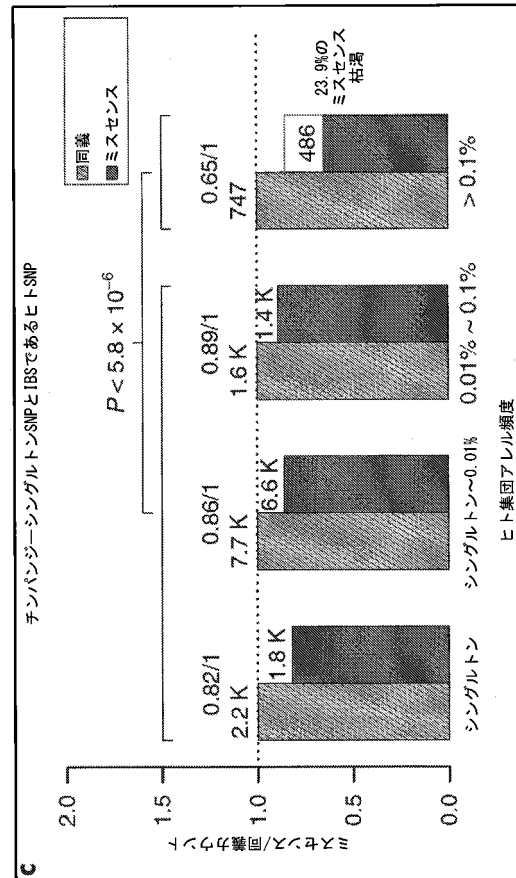
【 図 4 9 A 】



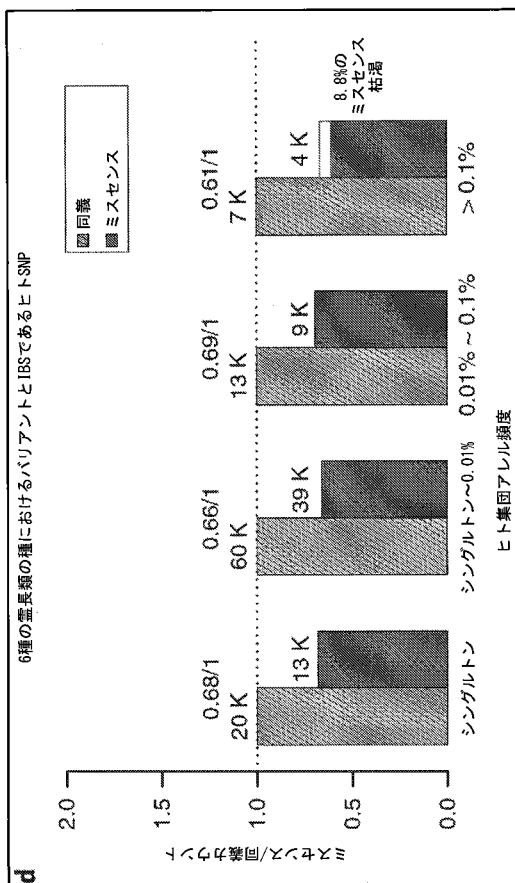
【図 49B】



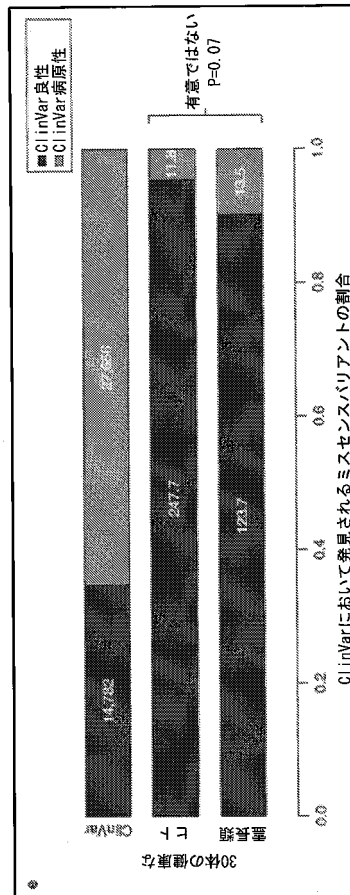
【図 49C】



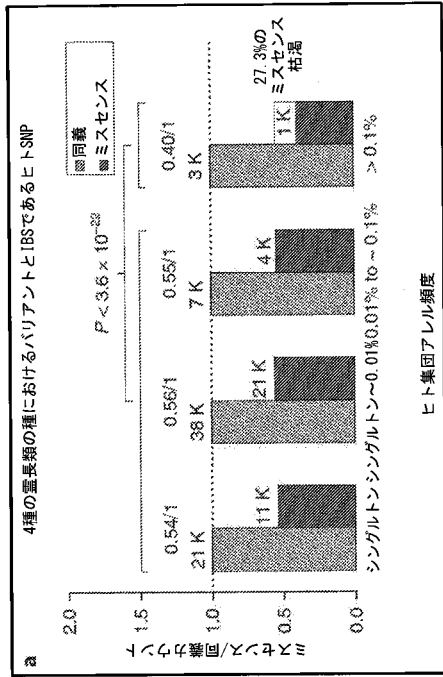
【図 49D】



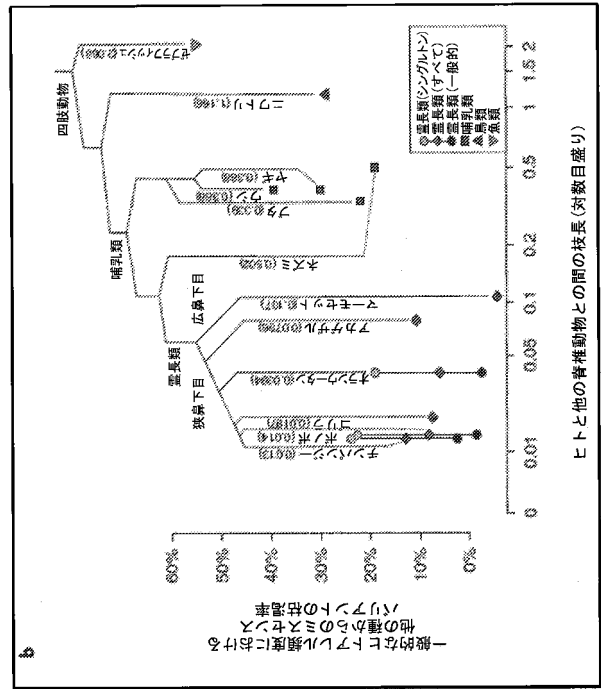
【図 49E】



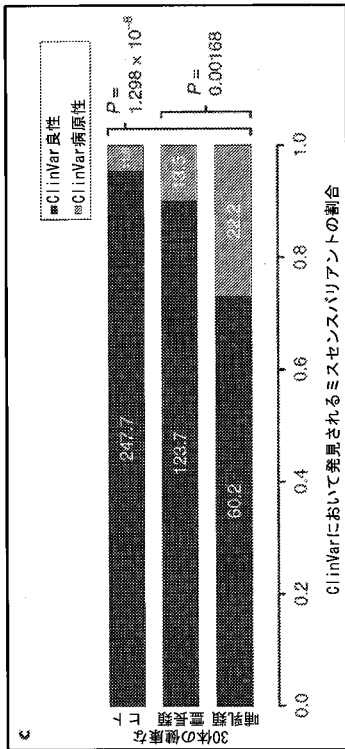
【図 50 A】



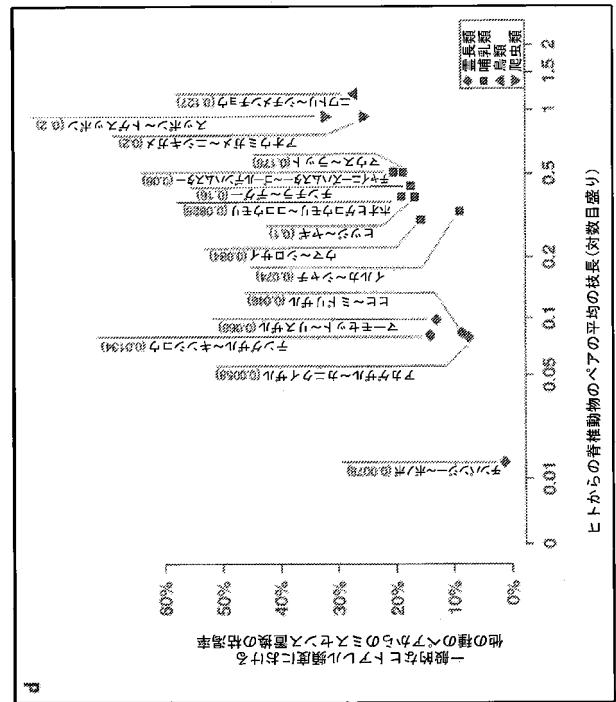
【図 50 B】



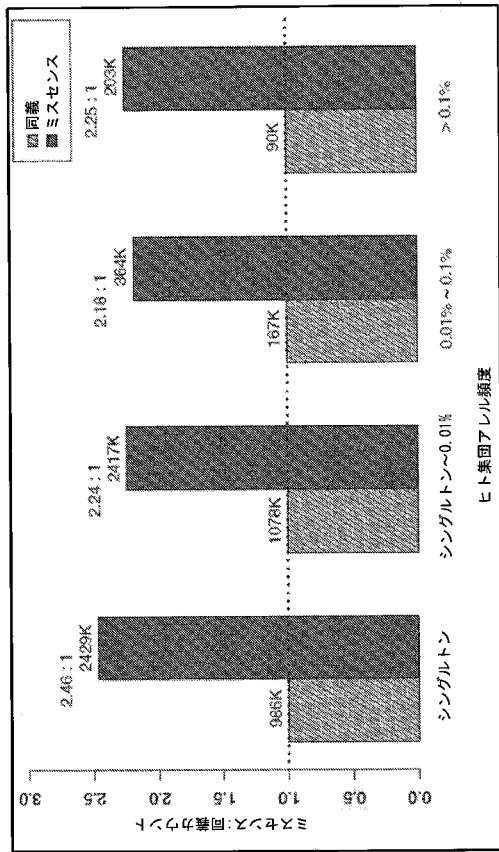
【図 50 C】



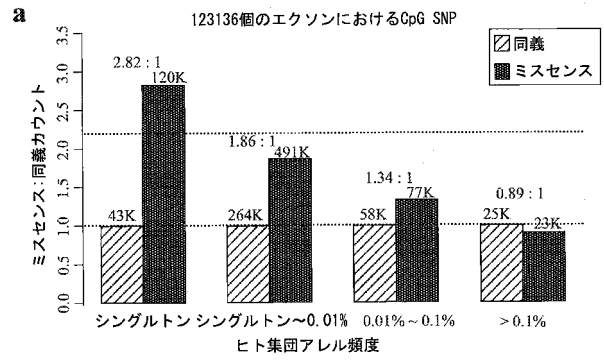
【図 50 D】



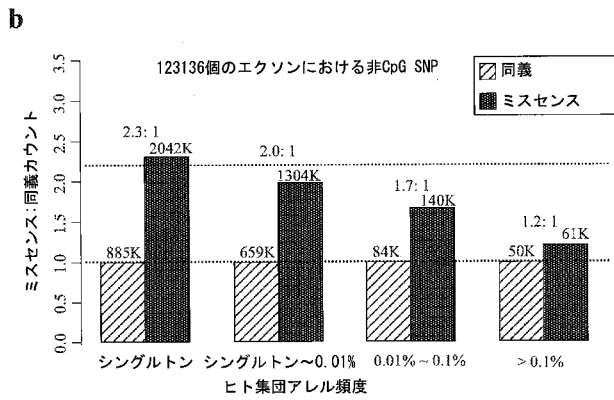
【図 5 1】



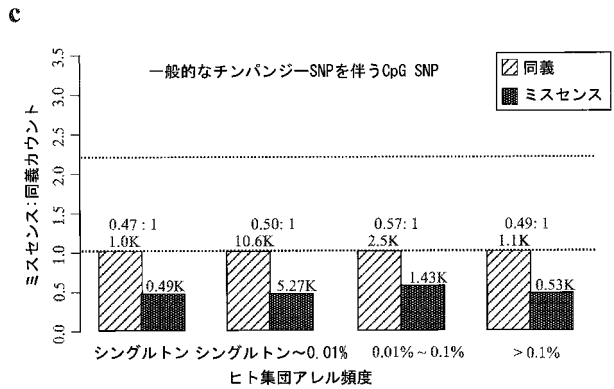
【図 5 2 A】



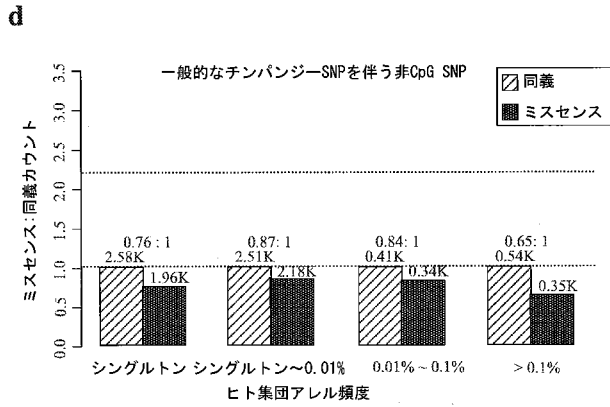
【図 5 2 B】



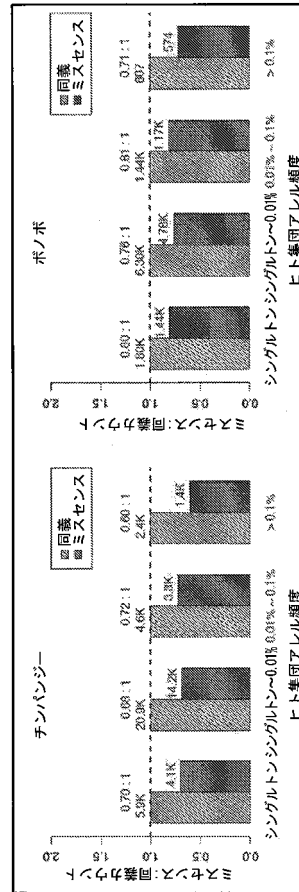
【図 5 2 C】



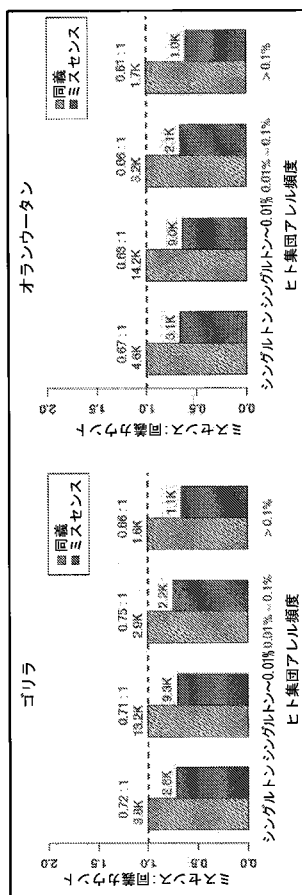
【 図 5 2 D 】



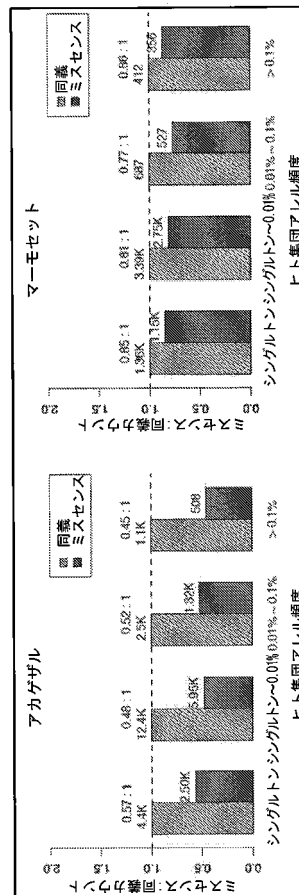
【 図 5 3 】



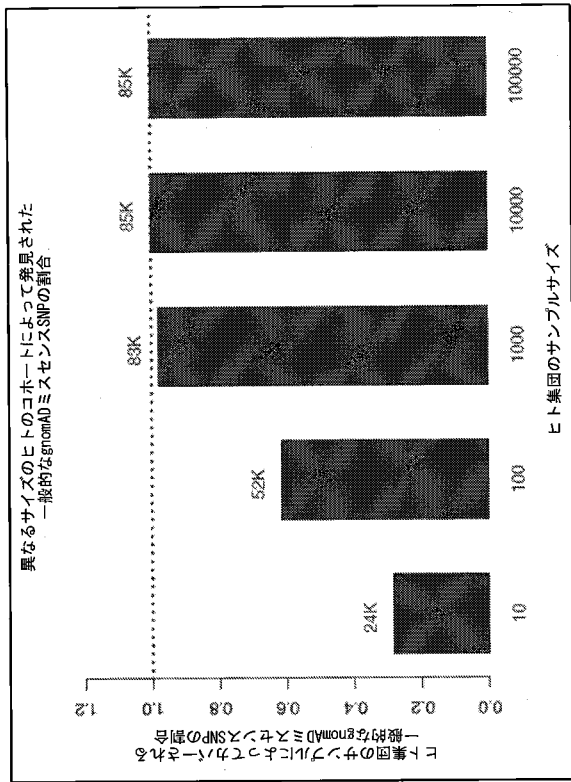
【 図 5 4 】



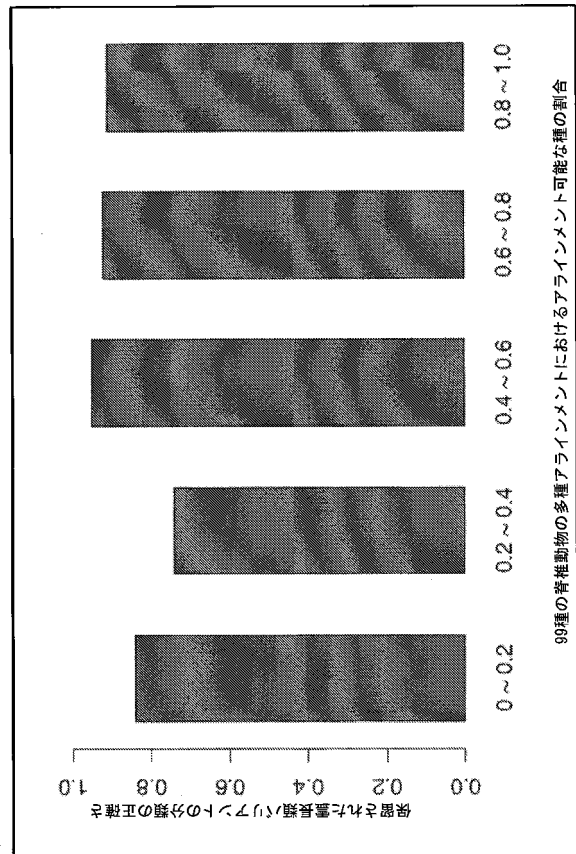
【 図 5 5 】



【図 5 6】



【図 5 7】



【図 5 8】

種 (および個体の数)	源	ヒトコドンと一致するミスセンスバリエーションの数	寄与する固有のミスセンスバリエーションの数
一般的なヒトバリエーション (MAF>0.1%、123136人の個人)	gnomAD/ExACデータベース	83,546	83,546
チンパンジー (24)	Prado-Martinez他、Nature (2013年)	76,293	73,858
ボノボ (13)	Prado-Martinez他、Nature (2013年)	25,277	22,968
ゴリラ (27)	Prado-Martinez他、Nature (2013年)	52,052	48,310
オランウータン、ボルネオ亜種 (5)	Prado-Martinez他、Nature (2013年)	21,621	19,793
オランウータン、スマトラ亜種 (5)	Prado-Martinez他、Nature (2013年)	27,567	19,834
チンパンジー (35)	de Manuel他、Science (2016年)	75,580	30,327
オランウータン (個体データ利用不可能)	dbSNP	18,627	10,554
アカゲザル (16、個体データ利用不可能)	dbSNP	62,393	53,279
マーモセット (9、個体データ利用不可能)	dbSNP	25,048	22,767
合計			385,236

【図 5 9】

123136個のエクソンにおけるヒト集団アレル頻度	同義バリエーションの予想される数	ミスセンスバリエーションの予想される数	ミスセンス/同義比
シングルトーン	986,141.22	2,429,287.85	2.46
シングルトーン~0.01%	1,077,630.37	2,416,751.61	2.24
0.01%~0.1%	166,849.26	364,292.75	2.18
>0.1%	90,221.42	202,918.38	2.24

【 図 6 0 】

元のClinVarファイイル	対象のClinVarバリアントの総数
hg19記録のみを保持する	666,403
一ヌクレオチドバリアントのみを抽出する	324,698
非コーディングバリアントを除去する	264,623
すべての同義変換と終止コドンを作成または破棄するミスセンス変換を無視する	186,769
有意性が知られていないバリアントおよび矛盾するアノテーションを伴うバリアントを除外する	122,884
	42,438

【 図 6 2 】

HGNシボル	タンパク質切断バリアント	ミスセンス	PrimateAIスコア $\geq 0.803$	すべてのミスセンス	PrimateAIスコア $\geq 0.803$	すべてのミスセンス	複数の個人において観察される表型型の異常
ACT168	0	3	1.5x10 <sup>-7</sup>	3	2.4x10 <sup>-4</sup>	3	小頭症
EBF3	3	3	5.2x10 <sup>-3</sup>	4	5.4x10 <sup>-3</sup>	4	生首遅延、眼奇形、斜視、運動失調
EFTLD2	2	4	1.5x10 <sup>-2</sup>	3	1.5x10 <sup>-2</sup>	3	小頭症、耳介低位、小耳症、後鼻孔閉鎖
HCCW2	1	6	2.3x10 <sup>-6</sup>	5	6.7x10 <sup>-7</sup>	5	発作、ミオパチー、頭蓋腔の異常
KDM6A	2	3	2.3x10 <sup>-6</sup>	3	9.8x10 <sup>-6</sup>	3	まぶた、歯牙の異常、筋緊張低下
RUSC	0	3	3.0x10 <sup>-4</sup>	3	2.5x10 <sup>-4</sup>	3	知的発達不全
RAP2K	0	5	3.1x10 <sup>-4</sup>	5	2.7x10 <sup>-4</sup>	5	面眼腫、耳介低位、羊水過多
PPP1C	0	6	1.5x10 <sup>-8</sup>	6	1.6x10 <sup>-8</sup>	6	顔の異常、低身長
PRKDI	0	6	8.6x10 <sup>-9</sup>	6	1.7x10 <sup>-8</sup>	6	肌、指、および心臓の異常、薄毛
SOX11	1	3	3.1x10 <sup>-4</sup>	3	2.4x10 <sup>-4</sup>	3	遺精、爪発育不全
TBR1	4	4	1.5x10 <sup>-8</sup>	4	4.2x10 <sup>-3</sup>	4	自閉症的行動
TK2	3	5	4.7x10 <sup>-3</sup>	5	6.3x10 <sup>-3</sup>	5	鼻、まぶたの異常、傾斜した聴覚線
TRIP2	5	2	1.4x10 <sup>-3</sup>	2	5.4x10 <sup>-3</sup>	2	関節弛緩
U2AF2	0	4	2.6x10 <sup>-3</sup>	4	3.2x10 <sup>-3</sup>	4	発作、眼、口蓋、人中の異常

表1| PrimateAIスコアが0.803以上であるミスセンスde novo変異(DNM)だけを考慮するときに知的障害者においてゲノムワイド有意性を達成する追加の遺伝子

タンパク質切断DNMおよびミスセンスDNMのカウントが提供される。PrimateAIスコアが0.803以上であるミスセンス変異だけをを用いて統計的検定が行われたとき、およびそれがすべてのミスセンス変異について繰り返されたときの、遺伝子エンリッチメントに対する計画が示される。

【 図 6 1 】

種	ClinVarにおける良性アノテーションを伴うバリアントの数	ClinVarにおける病原性アノテーションを伴うバリアントの数
チンパンジー	218	21
ボノボ	85	7
ゴリラ	167	23
オランウータン	160	12
アカゲザル	87	11
マーモセット	25	7
ネズミ	30	4
フタウシ	74	30
ウシ	77	39
ヤギ	60	16
ニワトリ	9	8
ゼブラフィッシュ	2	6

【 図 6 3 】

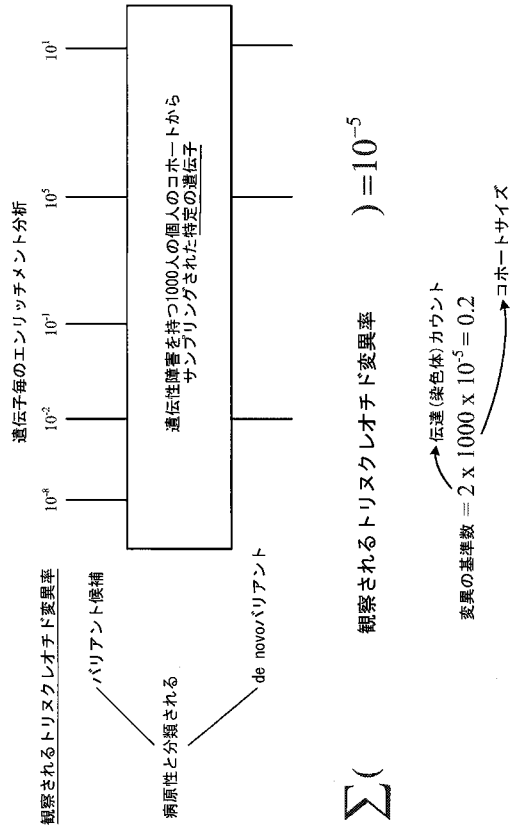
	Granthamスコア	タンパク質表面露出率	配列保存率
ClinVar病原性バリアント	911	0.53	0.87
ClinVar良性バリアント	674	0.41	0.54
専門家のアノテーションにおける差	+23.7	+0.12	+0.33
DDD患者におけるde novoバリアント	84.9	0.51	0.90
健康な対照群におけるde novoバリアント	72.7	0.29	0.73
影響を受けている個人と影響を受けていない個人の差	+12.2	+0.22	+0.17

表2| ClinVarにおける専門家によりアノテートされたバリアントと、DDD症例群vs対照群におけるde novoバリアントとの間での、Granthamスコア、タンパク質表面露出率、およびアミノ酸配列保存率の差の比較

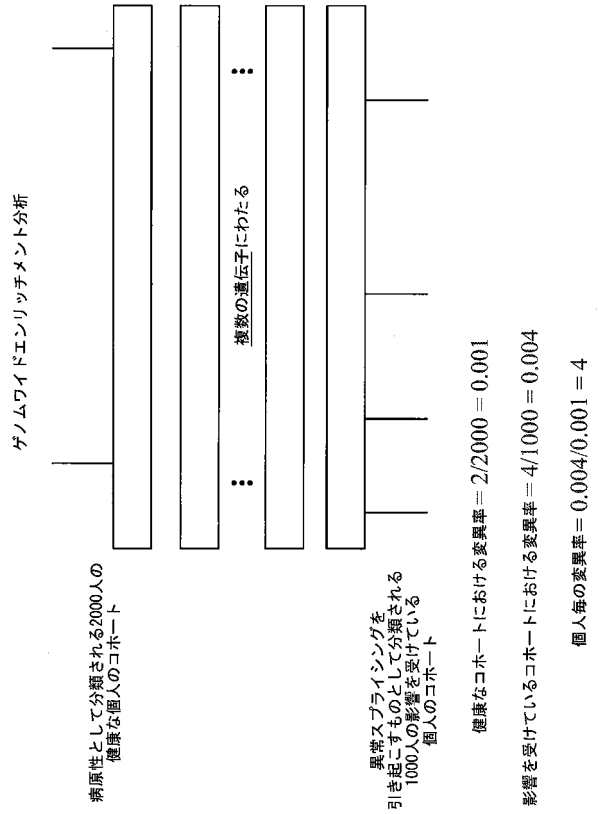
ClinVarデータベースの中の矛盾しないアノテーションを伴うミスセンス変異、および605個の疾病関連遺伝子内でDDD症例群vs対照群において存在するde novoバリアントに対する、平均スコアが示される。タンパク質表面露出率は、溶媒接触性ニューラルネットワークによって露出している残基として予測されるアミノ酸の割合を反映し、配列保存率は、100種の脊椎動物のアラインメントにおいて配列相同性を有するアミノ酸の割合を示す。太字の数字は、経験的データと比較して、専門家により好まれている経験則における差を強調する。



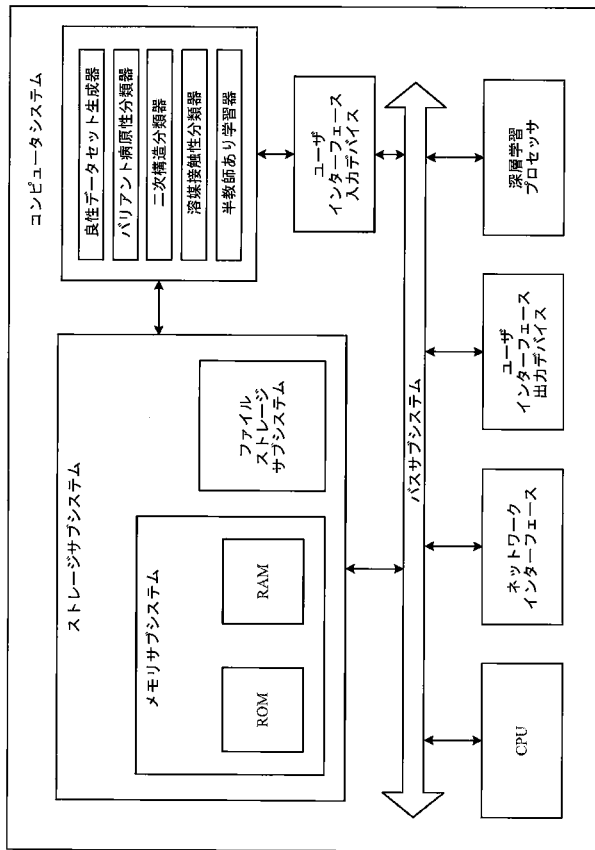
【 図 6 4 】



【 図 6 5 】



【 図 6 6 】



## 【 国際調査報告 】

## INTERNATIONAL SEARCH REPORT

International application No PCT/US2018/055840
---

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> INV. G16B20/20 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) G06F		
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,P	SUNDARAM LAKSSHMAN ET AL: "Predicting the clinical impact of human mutation with deep neural networks", NATURE GENETICS, NATURE PUBLISHING GROUP, NEW YORK, US, vol. 50, no. 8, 23 July 2018 (2018-07-23), pages 1161-1170, XP036657698, ISSN: 1061-4036, DOI: 10.1038/S41588-018-0167-Z [retrieved on 2018-07-23] abstract page 1161 - page 1168 ----- -/--	1-23
<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 "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 17 January 2019		Date of mailing of the international search report 25/01/2019
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer Kürten, Ivayla

2

## INTERNATIONAL SEARCH REPORT

International application No

PCT/US2018/055840

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	QIONG WEI ET AL: "The Role of Balanced Training and Testing Data Sets for Binary Classifiers in Bioinformatics", PLOS ONE, vol. 8, no. 7, 9 July 2013 (2013-07-09), page e67863, XP055541955, DOI: 10.1371/journal.pone.0067863	1-5, 15-23
A	abstract page 1 - page 4	6-14
Y	CA 2 894 317 A1 (DEEP GENOMICS INCORPORATED [CA]) 15 December 2016 (2016-12-15)	1-5, 15-23
A	abstract paragraph [0032] - paragraph [0036] paragraph [0052] - paragraph [0057] paragraph [0119] - paragraph [0126]	6-14
A	IOANNIDIS NILAH M ET AL: "REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants", AMERICAN JOURNAL OF HUMAN GENETICS, AMERICAN SOCIETY OF HUMAN GENETICS, CHICAGO, IL, US, vol. 99, no. 4, 22 September 2016 (2016-09-22), pages 877-885, XP029761034, ISSN: 0002-9297, DOI: 10.1016/J.AJHG.2016.08.016	1-23
A	abstract page 1581 - page 1582	
A	KARTHIK A JAGADEESH ET AL: "M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity", NATURE GENETICS., vol. 48, no. 12, 1 December 2016 (2016-12-01), pages 1581-1586, XP055541948, NEW YORK, US, ISSN: 1061-4036, DOI: 10.1038/ng.3703	1-23
	abstract page 878	

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No

PCT/US2018/055840

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
CA 2894317	A1	NONE	
-----			

## フロントページの続き

(31)優先権主張番号 62/573,153

(32)優先日 平成29年10月16日(2017.10.16)

(33)優先権主張国・地域又は機関  
米国(US)

(31)優先権主張番号 62/582,898

(32)優先日 平成29年11月7日(2017.11.7)

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

(72)発明者 カイ - ハウ・ファー

アメリカ合衆国・カリフォルニア・92122・サン・ディエゴ・イルミナ・ウェイ・5200

(72)発明者 ラクシュマン・サンダラム

アメリカ合衆国・カリフォルニア・92122・サン・ディエゴ・イルミナ・ウェイ・5200

(72)発明者 ジェレミー・フランシス・マクレー

アメリカ合衆国・カリフォルニア・92122・サン・ディエゴ・イルミナ・ウェイ・5200