

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2024-513569

(P2024-513569A)

(43)公表日 令和6年3月26日(2024.3.26)

(51)国際特許分類

G 0 6 F 40/44 (2020.01)

F I

G 0 6 F 40/44

審査請求 未請求 予備審査請求 未請求 (全29頁)

(21)出願番号	特願2023-561710(P2023-561710)	(71)出願人	312016539 ビットディフェンダー アイピーアール マネジメント リミテッド
(86)(22)出願日	令和4年3月28日(2022.3.28)		キプロス共和国、ニコシア、2012、 アクロポレオス、59-61、サード・ フロア、フラット/オフィス 302
(85)翻訳文提出日	令和5年11月30日(2023.11.30)	(74)代理人	100118902 弁理士 山本 修
(86)国際出願番号	PCT/EP2022/058130	(74)代理人	100106208 弁理士 宮前 徹
(87)国際公開番号	WO2022/214348	(74)代理人	100196508 弁理士 松尾 淳一
(87)国際公開日	令和4年10月13日(2022.10.13)	(72)発明者	アンドレイ・エム、マノラチェ ルーマニア国 615200 トゥルグ・ ネアムツ、ジュデトゥル・ネアムツ、ス 最終頁に続く
(31)優先権主張番号	17/301,641		
(32)優先日	令和3年4月9日(2021.4.9)		
(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,A T,BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR ,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC, 最終頁に続く		

(54)【発明の名称】 異常検出システムおよび方法

(57)【要約】

いくつかの実施形態が、自然言語処理やコンピュータセキュリティなどの適用分野での異常検出のために人工知能システム(たとえば、ディープニューラルネットワークのセット)をトレーニングする新規な手順を利用する。トレーニングコーパスから選択されたトークンシーケンスは、シーケンスアナライザに供給される前に、複数の所定のシーケンス変換のうち少なくとも1つに従ってひずめられる。さらに、シーケンスアナライザは、それぞれの入力トークンシーケンスを生成するためにどの変換が使用されたかを正しく推測するようにトレーニングされる。

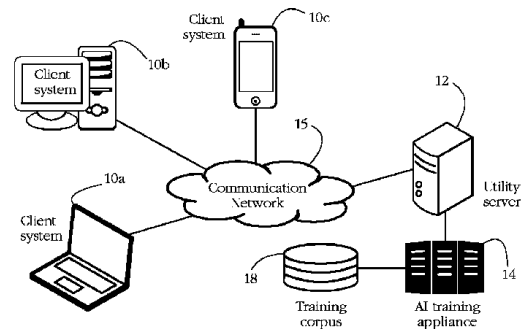


FIG. 1

**【特許請求の範囲】****【請求項 1】**

コンピュータで実装された異常検出方法であって、前記方法は、コンピュータシステムの少なくとも1つのハードウェアプロセッサを利用して、

トークンシーケンスのトレーニングコーパスからトレーニング・トークン・シーケンスを選択することに応答して、および、所定の複数のシーケンス変換から変換を選択することに応答して、選択された前記変換を前記トレーニング・トークン・シーケンスに適用して、修正後トークンシーケンスを生成するステップと、

調節可能パラメータのセットを有し、前記修正後トークンシーケンスに従って変換予測標識を決定するように構成されたシーケンスアナライザを実行するステップであって、前記変換予測標識は、選択された前記変換が適用されて前記修正後トークンシーケンスが生成された可能性を示す、シーケンスアナライザを実行するステップと、

前記予測標識を決定することに応答して、前記変換予測標識に従って調節可能パラメータの前記セットのうち少なくとも1つのパラメータを調節するステップと、

前記少なくとも1つのパラメータを調節することに応答して、前記シーケンスアナライザを利用してターゲット・トークン・シーケンスが異常であるかどうかを判定するステップと

を実施することを含む、コンピュータで実装された異常検出方法。

**【請求項 2】**

請求項 1 に記載の方法であって、選択された前記変換を適用するステップは、前記トレーニング・トークン・シーケンスの選択されたトークンを代替トークンと置き換えるステップを含む、方法。

**【請求項 3】**

請求項 2 に記載の方法であって、

調節可能パラメータの別のセットを有し、前記トレーニング・トークン・シーケンスに従って前記代替トークンを生成するように構成されたトークンジェネレータを実行するステップと、

前記変換予測標識を決定することに応答して、前記変換予測標識に従って調節可能パラメータの前記別のセットの別のパラメータを調節するステップと

をさらに含む方法。

**【請求項 4】**

請求項 1 に記載の方法であって、

選択された前記変換を適用するステップは、前記トレーニング・トークン・シーケンスの選択されたトークンを削除するステップと、前記トレーニング・トークン・シーケンス内に追加のトークンを挿入するステップと、前記トレーニング・トークン・シーケンスのトークンの選択されたサブセットを並べ替えるステップとを含むグループから選択された項目を含む、方法。

**【請求項 5】**

請求項 1 に記載の方法であって、

前記シーケンスアナライザは、前記修正後トークンシーケンスに従ってトークン予測標識を決定するようにさらに構成され、前記トークン予測標識は、前記修正後トークンシーケンスの選択されたトークンが、選択された前記変換の前記適用によって変更された可能性を示し、

前記少なくとも1つの調節可能パラメータを調節するステップは、前記トークン予測標識にさらに従って前記少なくとも1つの調節可能パラメータを調節するステップを含む、方法。

**【請求項 6】**

請求項 1 に記載の方法であって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、自然言語で構築されたテキストを含み、

10

20

30

40

50

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに対応して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスの作成者が前記トレーニング・トークン・シーケンスの作成者とは異なると判定するステップをさらに含む、方法。

【請求項 7】

請求項 1 に記載の方法であって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、自然言語で構築されたテキストを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに対応して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスの主題が前記トレーニング・トークン・シーケンスの主題とは異なると判定するステップをさらに含む、

方法。

【請求項 8】

請求項 1 に記載の方法であって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、自然言語で構築されたテキストを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに対応して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスがマシンで生成されたと判定するステップをさらに含む、

方法。

【請求項 9】

請求項 1 に記載の方法であって、

前記トレーニングコーパスは、選択基準に従って選択されたテキストフラグメントを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに対応して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスが前記選択基準を満たさないと判定するステップをさらに含む、

方法。

【請求項 10】

請求項 1 に記載の方法であって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、コンピューティングイベントのシーケンスを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに対応して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスがコンピュータセキュリティ脅威を示すと判定するステップをさらに含む、

方法。

【請求項 11】

少なくとも 1 つのハードウェアプロセッサを備えるコンピュータシステムであって、前記少なくとも 1 つのハードウェアプロセッサは、

トークンシーケンスのトレーニングコーパスからトレーニング・トークン・シーケンスを選択することに対応して、および、所定の複数のシーケンス変換から変換を選択することに対応して、選択された前記変換を前記トレーニング・トークン・シーケンスに適用して、修正後トークンシーケンスを生成することと、

調節可能パラメータのセットを有し、前記修正後トークンシーケンスに従って変換予測標識を決定するように構成されたシーケンスアナライザを実行することであって、前記変換予測標識は、選択された前記変換が適用されて前記修正後トークンシーケンスが生成さ

10

20

30

40

50

れた可能性を示す、シーケンスアナライザを実行することと、

前記予測標識を決定することに応答して、前記変換予測標識に従って調節可能パラメータの前記セットのうちの少なくとも1つのパラメータを調節することと、

前記少なくとも1つのパラメータを調節することに応答して、前記シーケンスアナライザを実行してターゲット・トークン・シーケンスが異常であるかどうかを判定することとを行うように構成された、コンピュータシステム。

【請求項12】

請求項11に記載のコンピュータシステムであって、選択された前記変換を適用することは、前記トレーニング・トークン・シーケンスの選択されたトークンを代替トークンと置き換えることを含む、コンピュータシステム。

10

【請求項13】

請求項12に記載のコンピュータシステムであって、前記少なくとも1つのハードウェアプロセッサは、

調節可能パラメータの別のセットを有し、前記トレーニング・トークン・シーケンスに従って前記代替トークンを生成するように構成されたトークンジェネレータを実行することと、

前記変換予測標識を決定することに応答して、前記変換予測標識に従って調節可能パラメータの前記別のセットの別のパラメータを調節することとを行うようにさらに構成される、コンピュータシステム。

【請求項14】

20

請求項11に記載のコンピュータシステムであって、

選択された前記変換を適用することは、前記トレーニング・トークン・シーケンスの選択されたトークンを削除することと、前記トレーニング・トークン・シーケンス内に追加のトークンを挿入することと、前記トレーニング・トークン・シーケンスのトークンの選択されたサブセットを並べ替えることとを含むグループから選択された項目を含む、コンピュータシステム。

【請求項15】

請求項11に記載のコンピュータシステムであって、

前記シーケンスアナライザは、前記修正後トークンシーケンスに従ってトークン予測標識を決定するようにさらに構成され、前記トークン予測標識は、前記修正後トークンシーケンスの選択されたトークンが選択された前記変換の前記適用によって変更された可能性を示し、

30

前記少なくとも1つの調節可能パラメータを調節することは、前記トークン予測標識に従って前記少なくとも1つの調節可能パラメータを調節することを含む、コンピュータシステム。

【請求項16】

請求項11に記載のコンピュータシステムであって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、自然言語で構築されたテキストを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに応答して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスの作成者が前記トレーニング・トークン・シーケンスの作成者とは異なると判定することをさらに含む、コンピュータシステム。

40

【請求項17】

請求項11に記載のコンピュータシステムであって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、自然言語で構築されたテキストを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに応答して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲ

50

ット・トークン・シーケンスの主題が前記トレーニング・トークン・シーケンスの主題とは異なると判定することをさらに含む、  
コンピュータシステム。

【請求項 18】

請求項 11 に記載のコンピュータシステムであって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、自然言語で構築されたテキストを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに応答して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスがマシンで生成されたと判定することをさらに含む、  
コンピュータシステム。

10

【請求項 19】

請求項 11 に記載のコンピュータシステムであって、

前記トレーニングコーパスは、選択基準に従って選択されたテキストフラグメントを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに応答して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスが前記選択基準を満たさないと判定することをさらに含む、  
コンピュータシステム。

20

【請求項 20】

請求項 11 に記載のコンピュータシステムであって、

前記トレーニング・トークン・シーケンスおよびターゲット・トークン・シーケンスは、コンピューティングイベントのシーケンスを含み、

前記方法は、前記ターゲット・トークン・シーケンスが異常であるかどうかを判定することに応答して、前記ターゲット・トークン・シーケンスが異常であるとき、前記ターゲット・トークン・シーケンスがコンピュータセキュリティ脅威を示すと判定することをさらに含む、  
コンピュータシステム。

【請求項 21】

命令を記憶する非一時的コンピュータ可読媒体であって、前記命令は、コンピュータシステムの少なくとも 1 つのハードウェアプロセッサによって実行されるとき、前記コンピュータシステムに、

30

トークンシーケンスのトレーニングコーパスからトレーニング・トークン・シーケンスを選択することに応答して、および、所定の複数のシーケンス変換から変換を選択することに応答して、選択された前記変換を前記トレーニング・トークン・シーケンスに適用して、修正後トークンシーケンスを生成することと、

調節可能パラメータのセットを有し、前記修正後トークンシーケンスに従って変換予測標識を決定するように構成されたシーケンスアナライザを実行することであって、前記変換予測標識は、選択された前記変換が適用されて前記修正後トークンシーケンスが生成された可能性を示す、シーケンスアナライザを実行することと、

40

前記予測標識を決定することに応答して、前記変換予測標識に従って調節可能パラメータの前記セットのうちの少なくとも 1 つのパラメータを調節することと、

前記少なくとも 1 つのパラメータを調節することに応答して、前記シーケンスアナライザを実行してターゲット・トークン・シーケンスが異常であるかどうかを判定することとを行わせる、非一時的コンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

[0001]本発明は人工知能に関し、詳細には自然言語処理およびコンピュータセキュリティの適用分野のための、データ内の異常を自動的に検出するためのシステムおよび方法

50

に関する。

【背景技術】

【0002】

【0002】人工知能（AI）技術および機械学習技術は、とりわけパターン認識、自動分類、および異常検出などの適用分野に関して大量のデータを処理するためにますます使用されている。異常検出は、基準グループによって集合的に定義された標準または「正常」から大幅に逸脱する標本を識別することに相当する。異常検出は、複雑なデータのケースではかなりの技術的課題を提起し得、そのケースでは、正常性の意味および境界が事前に明確ではないことがあり、または事前に定義されないことがある。データから精巧なモデルを自動的に推論する能力を用いて、現代の人工知能システム（たとえばディープニューラルネットワーク）は、そのような課題に対して良好に動作することが示されている。

10

【0003】

【0003】しかしながら、異常検出器をトレーニングするために機械学習を実装することは、それ自体の技術的課題のセットを提起する。従来手法のいくつかでは、トレーニングが極度の計算コストを招くことがあり、非常に大規模なトレーニングコーパスを必要とすることがあり、不安定および/または非効率であることがある。したがって、自然言語処理およびコンピュータセキュリティの適用分野について異常検出器をトレーニングする新規な検出器アーキテクチャおよび新規な方法を開発することに関心向けられている。

20

【発明の概要】

【0004】

【0004】一態様によれば、コンピュータで実装された異常検出方法が、トークンシーケンスのトレーニングコーパスからトレーニング・トークン・シーケンスを選択することに応答して、かつ所定の複数のシーケンス変換から変換を選択することに応答して、コンピュータシステムの少なくとも1つのハードウェアプロセッサを利用して、選択された変換をトレーニング・トークン・シーケンスに適用して、修正後トークンシーケンスを生成するステップを含む。方法は、調節可能パラメータのセットを有し、修正後トークンシーケンスに従って変換予測標識を決定するように構成されたシーケンスアナライザを実行するステップであって、変換予測標識は、選択された変換が適用されて修正後トークンシーケンスが生成された可能性を示す、シーケンスアナライザを実行するステップをさらに含む。

方法は、予測標識を決定することに応答して、変換予測標識に従って調節可能パラメータのセットのうち少なくとも1つのパラメータを調節するステップと、少なくとも1つのパラメータを調節することに応答して、シーケンスアナライザを利用してターゲット・トークン・シーケンスが異常であるかどうかを判定するステップとをさらに含む。

30

【0005】

【0005】別の態様によれば、コンピュータシステムが、トークンシーケンスのトレーニングコーパスからトレーニング・トークン・シーケンスを選択することに応答して、かつ所定の複数のシーケンス変換から変換を選択することに応答して、選択された変換をトレーニング・トークン・シーケンスに適用して、修正後トークンシーケンスを生成することを行うように構成された少なくとも1つのハードウェアプロセッサを備える。少なくとも1つのハードウェアプロセッサは、調節可能パラメータのセットを有し、修正後トークンシーケンスに従って変換予測標識を決定するように構成されたシーケンスアナライザを実行することであって、変換予測標識は、選択された変換が適用されて修正後トークンシーケンスが生成された可能性を示す、シーケンスアナライザを実行することを行うようにさらに構成される。少なくとも1つのハードウェアプロセッサは、予測標識を決定することに応答して、変換予測標識に従って調節可能パラメータのセットのうち少なくとも1つのパラメータを調節することと、少なくとも1つのパラメータを調節することに応答して、シーケンスアナライザを利用してターゲット・トークン・シーケンスが異常であるかどうかを判定することを行うようにさらに構成される。

40

【0006】

50

【0006】別の態様によれば、非一時的コンピュータ可読媒体が、コンピュータシステムの少なくとも1つのハードウェアプロセッサによって実行されるとき、コンピュータシステムに、トークンシーケンスのトレーニングコーパスからトレーニング・トークン・シーケンスを選択することに応答して、かつ所定の複数のシーケンス変換から変換を選択することに応答して、選択された変換をトレーニング・トークン・シーケンスに適用して、修正後トークンシーケンスを生成させることを行わせる命令を記憶する。命令はさらに、コンピュータシステムに、調節可能パラメータのセットを有し、修正後トークンシーケンスに従って変換予測標識を決定するように構成されたシーケンスアナライザを実行することによって、変換予測標識は、選択された変換が適用されて修正後トークンシーケンスが生成された可能性を示す、シーケンスアナライザを実行することを行わせる。命令はさらに、コンピュータシステムに、予測標識を決定することに応答して、変換予測標識に従って調節可能パラメータのセットのうち少なくとも1つのパラメータを調節させることと、少なくとも1つのパラメータを調節することに応答して、シーケンスアナライザを利用してターゲット・トークン・シーケンスが異常であるかどうかを判定することを行わせる。

10

#### 【0007】

【0007】以下の詳細な説明を読み、図面を参照するとき、本発明の前述の態様および利点をより良く理解されよう。

#### 【図面の簡単な説明】

#### 【0008】

【図1】【0008】本発明のいくつかの実施形態による、異常を検出する際にユーティリティサーバと協働するクライアントシステムのセットを示す図である。

20

【図2】【0009】本発明のいくつかの実施形態による、異常検出器の例示的動作を示す図である。

【図3】【0010】本発明のいくつかの実施形態による、異常検出器の例示的トレーニングを示す図である。

【図4】【0011】本発明のいくつかの実施形態による、入力修正器の例示的動作を示す図である。

【図5】【0012】本発明のいくつかの実施形態による例示的トークン埋込み空間を示す図である。

30

【図6】【0013】いくつかの実施形態による例示的シーケンス変換を示し、図示される変換が、選択されたトークンの代表的ベクトルをナッジすることを含む図である。

【図7】【0014】本発明のいくつかの実施形態によるシーケンス分類器の例示的構造を示す図である。

【図8】【0015】本発明のいくつかの実施形態による、異常検出器のトレーニング中に実施されるステップの例示的シーケンスを示す図である。

【図9】【0016】本発明のいくつかの実施形態による、トレーニング済み異常検出器によって実施されるステップの例示的シーケンスを示す図である。

【図10】【0017】本明細書で説明される方法のうちいくつかを実施するように構成された例示的コンピューティングアプライアンスを示す図である。

40

#### 【発明を実施するための形態】

#### 【0009】

【0018】以下の説明では、構造間のすべての記載の接続が直接的に動作可能な接続または中間構造を通じた間接的に動作可能な接続であり得ることを理解されたい。要素のセットは1つまたは複数の要素を含む。要素のどんな記載も少なくとも1つの要素を指すと理解されたい。複数の要素は少なくとも2つの要素を含む。別段に指定されていない限り、「または(OR)」のどんな使用も非排他的論理和を指す。別段に必要とされない限り、記載のどんな方法ステップも、必ずしも特定の図示される順序で実施される必要はない。第2の要素から導出された第1の要素(たとえばデータ)は、第2の要素と等しい第1の要素、ならびに第2の要素と、任意選択で他のデータとを処理することによって生成され

50

た第1の要素を包含する。パラメータに従って決定または判断を行うことは、パラメータに従って、かつ任意選択で他のデータに従って決定または判断を行うことを包含する。別段に指定されていない限り、いくつかの量/データの標識は、量/データ自体、または量/データ自体とは異なる標識であり得る。コンピュータプログラムは、タスクを実施するプロセッサ命令のシーケンスである。本発明のいくつかの実施形態で説明されるコンピュータプログラムは、スタンドアロン・ソフトウェア・エンティティまたは他のコンピュータプログラムのサブエンティティ（たとえば、サブルーチン、ライブラリ）であり得る。コンピュータ可読媒体は、磁気、光、および半導体記憶媒体（たとえば、ハードドライブ、光ディスク、フラッシュメモリ、DRAM）などの非一時的媒体、ならびに導電性ケーブルや光ファイバリンクなどの通信リンクを包含する。いくつかの実施形態によれば、本発明は、とりわけ、本明細書に記載の方法を実施するようにプログラムされるハードウェア（たとえば、1つまたは複数のプロセッサ）、ならびに本明細書に記載の方法を実施するための命令を符号化するコンピュータ可読媒体を備えるコンピュータシステムを提供する。

10

#### 【0010】

[0019]以下の説明は、必ずしも限定としてではなく、例として本発明の実施形態を示す。

[0020]図1は、本発明のいくつかの実施形態による、データ内の異常を検出するためにユーティリティサーバ12と協働し得るクライアントシステム10a~cの例示的セットを示す。本明細書では、異常とは、項目の基準集合/コーパスによって集合的に表される標準または「正常」から大幅に逸脱する項目を表すと理解されたい。この説明は、異常なテキストフラグメントやコンピューティングイベントシーケンスなどの異常なトークンシーケンスを検出することに焦点を当てる。そのような実施形態では、例示的な異常検出は、ターゲットテキストの作成者が基準テキストとは異なると判定することを含む。別の例示的な異常検出は、コンピューティングイベントのシーケンスがそれぞれのコンピュータの正常な挙動から逸脱し、恐らくはセキュリティ違反または悪意のあるソフトウェアの存在を示すと判定することを含む。いくつかの例示的な異常検出使用事例シナリオが以下で説明される。

20

#### 【0011】

[0021]クライアントシステム10a~cは一般に、プロセッサ、メモリ、および通信インターフェースを有する任意の電子アプライアンスを表す。例示的クライアントシステム10a~cは、とりわけ、パーソナルコンピュータ、企業メインフレームコンピュータ、サーバ、ラップトップ、タブレットコンピュータ、モバイル遠隔通信デバイス（たとえば、スマートフォン）、メディアプレーヤ、TV、ゲームコンソール、ホームアプライアンス、およびウェアラブルデバイス（たとえば、スマートウォッチ）を含む。図示されるクライアントシステムは通信ネットワーク15によって相互接続され、通信ネットワーク15は、ローカルエリアネットワーク(LAN)および/またはインターネットなどの広域ネットワーク(WAN)を含み得る。サーバ12は一般に、互いに物理的に近接することがあり、または近接しないことがある、通信可能に結合されたコンピュータシステムのセットを表す。

30

40

#### 【0012】

[0022]図2は、本発明のいくつかの実施形態による例示的異常検出器20の動作を示す。異常検出器20は、ソフトウェアとして、すなわちメモリ内にロードされ、パーソナルコンピュータやスマートフォンなどのコンピューティングアプライアンスのハードウェアプロセッサによって実行されるとき、それぞれのアプライアンスにそれぞれのタスクを実施させる命令を含むコンピュータプログラムのセットとして実施され得る。しかしながら、そのような実施形態が限定を意味するわけではないことを当業者は理解されよう。むしろ、検出器20は、ソフトウェアとハードウェアの任意の組合せとして実装され得る。たとえば、検出器20の一部またはすべての機能が、フィールドプログラマブルゲートアレイ(FPGA)や他の特定用途向け集積回路(ASIC)などのファームウェアおよび

50



／または専用ハードウェアで実装され得る。それぞれのハードウェアモジュールは、それぞれの機能向けに高度に最適化され、たとえば、特定のバージョンのディープニューラルネットワークアーキテクチャを直接的に実装し、したがって汎用プロセッサに関して達成可能なものよりもかなり速い処理速度を可能にし得る。さらに、以下で説明されるように異常検出器 20 および／または検出器 20 をトレーニングするように構成されたコンピュータシステムの別個の構成要素が、別個であるが通信可能に結合されたマシン上で、かつ／または同一のコンピュータシステムの別個のハードウェアプロセッサ上で実行され得ることを当業者は理解されよう。

#### 【0013】

[0023]異常検出器 20 は、ターゲット・トークン・シーケンス 22 のコンピュータ可読符号化を受け取り、それに応答して、それぞれのトークンシーケンス 22 が異常であるかどうかを示す異常標識 26 を出力するように構成され得る。例示的トークンシーケンスは、とりわけ、英語や中国語などの自然言語で構築されたテキストのフラグメントなどのトークンの順序付き配列を含む。一般性を失うことなく、以下の説明は、主に自然言語処理の例に焦点を当て、例示的トークンは、とりわけ、個々の語、語句、文、数、句読点（たとえば、？！；：／（）, . . . ）、特殊文字（たとえば、\$ # % ）、省略語（USA、LOL、IMHO など）、ソーシャルメディアハンドル（たとえば、@POTUS ）、ハッシュタグ、エモティコンを含み得る。本明細書で説明されるシステムおよび方法は、とりわけ、コンピューティングイベントのシーケンスやサウンドシーケンス（たとえば、音楽、音声）などの他のタイプのトークンシーケンスを処理することに適合され得ることを当業者は理解されよう。

#### 【0014】

[0024]例示的異常標識 26 は、それぞれのターゲット・トークン・シーケンスが異常である可能性を示す数値スコアを含む。スコアはプール（たとえば、YES / NO）であり得、または所定の境界の間（たとえば、0 から 1 の間）で徐々に変動し得る。そのような一例では、大きい値は、それぞれのシーケンスが異常である可能性が高いことを示す。代替の異常標識 26 は、シーケンス 22 が属する可能性が高いトークンシーケンスのカテゴリを示す分類ラベル（たとえば、異常、正常、未知、容疑）を含み得る。

#### 【0015】

[0025]1つの例示的シナリオでは、異常検出器 20 の別個のインスタンスが、各クライアントシステム 10 a ~ c 上で実行され得、したがって各クライアントは、それ自体の異常検出活動をローカルに独立して実施し得る。代替実施形態では、異常検出器 20 はユーティリティサーバ 12 上で実行され得、したがってユーティリティサーバ 12 は、複数のクライアントシステム 10 a ~ c の代わりに集中型異常検出活動を実施し得る。そのような実施形態では、サーバ 12 は、各クライアントシステム 10 a ~ c からターゲット・トークン・シーケンス 22 の符号化を受信し、それぞれの異常標識 26 をそれぞれのクライアントに返し得る。そのような一例では、クライアント 10 a ~ c は、ユーティリティサーバ 12 によって公開されるウェブインターフェースを介して異常検出サービスにアクセスし得る。

#### 【0016】

[0026]図 3 は、判断モジュール 44 に接続されたシーケンスアナライザ 42 などの異常検出器の例示的構成要素を示す。いくつかの実施形態では、シーケンスアナライザ 42 は、基準トークンシーケンスのコーパス 18 に関してトレーニングされたディープニューラルネットワークなどの人工知能（AI）システムを備える。自然言語処理のシナリオでは、コーパス 18 は、自然言語（たとえば、英語）で書かれたテキストフラグメントの集合を含み得る。コーパス 18 のより具体的な例は、特定の作成者によるテキストの集合、電子メッセージ（たとえば、ショートメッセージサービス - SMS メッセージ、eメール、ソーシャルメディアポストなど）の集合、特定の話題または関心のエリア（たとえば、ビジネスニュース、スポーツ、中東など）に関するテキストの集合、および特定のスタイル（たとえば、フィクション、詩、科学記事、ニュースなど）で書かれたテキストの集合

10

20

30

40

50

からなり得る。個々のコーパス項目は、たとえばメタデータを使用して、タグ付けされ、ラベル付けされ、かつ/または注釈付けされ得る。例示的メタデータは、項目の選択されたクラス/カテゴリ(たとえば、特定のユーザによって送られたeメールメッセージ、金融ニュースなど)に対するメンバシップの標識を含み得る。コーパス18は、当技術分野で周知の任意のフォーマットで、たとえばリレーショナルデータベース、単純リスト、またはXMLもしくはJSONフォーマットで指定された構造化データとして編成され、記憶され得る。

#### 【0017】

[0027]コーパス18の内容は、通信の基準パターンまたは「正常な」パターンを集合的に定義し、いくつかの実施形態では、異常検出器20は、それぞれの基準パターンの内部モデルを構築し、それに応答して、ターゲット・テキスト・フラグメントが学習済みのパターンに適合するか否かを判定するようにトレーニングされ得る。ターゲット・トークン・シーケンス22が(コーパス18に従って)「正常な」テキストに対応する基準パターンに適合しないことが判明したとき、シーケンス22は異常であると見なされ、異常標識26を介してそのようなものとしてレポートされ得る。

10

#### 【0018】

[0028]いくつかの実施形態では、異常検出器20のトレーニングが、図1のAイトレーニングアプライアンス14として示される、別々の専用コンピュータシステムによって実施される。アプライアンス14は、ユーティリティサーバ12および/またはクライアントシステム10a~cに通信可能に結合され得、計算コストのかかるトレーニング手順を容易にするための、グラフィックス処理装置(GPU)ファームなどの専用ハードウェアを備え得る。「トレーニング」という用語は通常、当技術分野では機械学習手順を示すために使用され、それによって、人工知能システム(たとえば、ニューラルネットワーク)に様々なトレーニング入力提示され、それぞれの入力生成する出力に従って人工知能システムが徐々に調整される。各トレーニング入力/バッチについて、トレーニングは、それぞれの入力を処理してトレーニング出力を生成することと、それぞれのトレーニング出力および/または入力に従って問題特有のユーティリティ関数の値を求めることと、それぞれのユーティリティ値に従ってそれぞれのAIシステムのパラメータのセットを調節することとを含み得る。パラメータを調節することは、ユーティリティ関数を最大化する(いくつかのケースでは、最小化する)ことを目標とし得る。ニューラルネットワークをトレーニングする一例では、調節可能パラメータはシナプス重みのセットを含み得、一方、ユーティリティ関数は、予想される出力または所望の出力からのトレーニング出力の逸脱を定量化し得る。そのような一例では、トレーニングは、それぞれのトレーニング入力に対応する所望の出力にトレーニング出力を近づけるように、シナプス重み、および場合によっては他のネットワークパラメータを調節することを含み得る。既知のトレーニングのフレーバは、とりわけ、教師あり、教師なし、自己教師あり、および強化学習を含む。いくつかの実施形態では、典型的な検出器20の調節可能パラメータの数は、数千から数百万まで様々であり得る。トレーニングが成功すると、最適化された検出器パラメータ値24(図2)のセットを生成し得、最適化された検出器パラメータ値24は、クライアントシステム10a~cおよび/またはユーティリティサーバ12上で実行中の異常検出器20のローカルインスタンスをインスタンス化するために使用され得る。

20

30

40

#### 【0019】

[0029]検出器20のトレーニングが、図3に概略的に示されている。Aイトレーニングアプライアンス14の同一のハードウェアプロセッサまたは物理マシン上で実行するために、図示される構成要素のすべてが必要であるわけではないことを当業者は理解されよう。

#### 【0020】

[0030]本発明のいくつかの実施形態は、トレーニングコーパス18内に含まれるサンプルのうち少なくともいくつかを、異常検出器20内に供給する前にひずませ、次いで検出器20をトレーニングして、適用されたひずみのタイプを識別する。図3に示される

50

一例では、入力修正器40が、トレーニングコーパス18から選択されたトレーニング・トークン・シーケンス32を受け取り、トレーニングシーケンス32にシーケンス変換30の所定のセットのうちの少なくとも1つを適用した結果を含む修正後トークンシーケンス34を出力するように構成される。

#### 【0021】

[0031]例示的シーケンス変換30は、とりわけ、シーケンス32内のトークンの選択されたサブセットを代替トークンと置き換えることと、シーケンス32からトークンの選択されたサブセットを削除することと、トークンのセットをシーケンス32内に挿入することと、シーケンス32内のトークンの選択されたサブセットを並べ替えることとを含む。それぞれの交換による修正の目標とされるトークンのサブセットは、トレーニングシーケンス内の各トークンの位置に従って選択され得る。目標とされる位置はバイナリマスクによって示され得、0が、不変のままにされるトークンの位置をマークし、1が、それぞれの変換によって影響を受ける位置をマークする。たとえば、マスク[0 0 1 0 1]によって定義される並べ替え変換は、トークンシーケンス「They were prepared to leave」を修正後トークンシーケンス「They were leave to prepared」に変換し得、第3のトークンが第5のトークンと交換された。

#### 【0022】

[0032]代替実施形態では、変換30によって目標とされるトークンは、それぞれのトークンのタイプに従って選択され得る。たとえば、いくつかの変換は、特定の品詞（たとえば、名詞、動詞、形容詞）または特定の文法的役割（たとえば、文の主語）を有するトークンを目標とし得る。1つのそのような例示的変換は、動詞を代替動詞または動詞句で置き換え得る。それぞれの代替トークンまたはトークンシーケンスは、ターゲットトークン/シーケンスの同意語または反意語となるように選択され得る。シーケンス変換30のより高度な例は、パラフレーズング、すなわち意味を保持しながらトークンシーケンス全体を代替シーケンスと置き換えることを含み得る。パラフレーズング変換の一例は、トレーニングシーケンス「Kids by the lake were being eaten alive by mosquitoes」を修正後シーケンス「Mosquitoes ferociously attacked the kids by the lake」と置き換えることを含む。

#### 【0023】

[0033]同様の交換が、シーケンス32の各トークンが個々のコンピューティングイベントを含む、コンピュータセキュリティ実施形態で適用され得る。たとえば、例示的変換30が、タイプ「create process」のトークンをトレーニングシーケンス32から除去し得る。そのような実施形態でのパラフレーズングの同等物が、イベントのターゲットシーケンスを、それぞれのコンピュータシステムを同一の最終状態にすることになるイベントの代替シーケンスと置き換えることを含み得る。

#### 【0024】

[0034]図4は、本発明のいくつかの実施形態による、複数の事前定義されたシーケンス変換30を実装する入力修正器40の例示的動作を示す。トレーニングシーケンス32はトークン35a~eを含み、現在の例では、トークン35a~eは個々の語である。図示される実施形態では、各変換30は、トレーニングシーケンス32からのトークンのセットを代替トークンと置き換えることを含む。各変換30は、トークン置換えの目標とされるトレーニングシーケンス32内の位置の別個のセットを示す、図示されるような別個のバイナリマスクによって定義され得る。図示される例では、T2マスクをシーケンス32に適用することは、シーケンス32の第3および第5のトークンをマスクし、実質的に前記トークンを置換えのためにマークする。

#### 【0025】

[0035]いくつかの実施形態では、入力修正器40は、代替トークンのセットを出力して、シーケンス32内のマスクされたトークンを置き換えるように構成されたトークンジ

10

20

30

40

50

ジェネレータ 4 1 をさらに備える。図示される例では、トークンジェネレータは、代替トークン 3 5 f および 3 5 g を出力して、それぞれトークン 3 5 c および 3 5 e を置き換える。ジェネレータ 4 1 の単純な実施形態は、基準プールから代替トークンをランダムに引き出すように構成され得る。より高度な実施形態では、トークンジェネレータ 4 1 は辞書 / シソーラスを備え、それぞれのマスクされたトークンについて、それぞれのトークンの同意語または反意語を出力するように構成され得る。別の例示的实施形態では、ジェネレータ 4 1 は、マスクされたトークンに従って、さらにはマスクされたトークンのコンテキストに従って代替トークンを決定し得、コンテキストは、マスクされたトークンに先行するシーケンス 3 2 のフラグメント、および / またはマスクされたトークンの後のシーケンス 3 2 のフラグメントからなる。たとえば、図 4 の例では、トークン 3 5 c (「ready」) のコンテキストは、トークン 3 5 b および 3 5 d (「were to」) を含み得る。そのような実施形態は、統計言語モデルを利用して、マスクされたトークンのコンテキスト内のそれぞれの代替トークンの発生の確率に従って代替トークンを生成し得る。言い換えれば、トークンジェネレータ 4 1 は、それぞれのトークンシーケンスのコンテキストが与えられると、妥当な代替トークンを生成するように構成され得る。

#### 【0026】

[0036] 妥当な代替トークンを生成するトークンジェネレータ 4 1 の例示的一実施形態は、それぞれの異常検出アプリケーションを表すトークンシーケンスのコーパスに関してトレーニングされた AI システム (たとえば、ディープニューラルネットワークのセット) を備える。ジェネレータ 4 1 のそのようなバージョンは、シーケンス 3 2 内のマスクされたトークンに先行するトークンのサブシーケンスに従って代替トークンを出力し得る。図示される例では、トレーニング済みジェネレータ 4 1 は、トークン 3 5 a ~ d のシーケンス (「they were ready to」) に続く高い可能性を有するものとして代替トークン 3 5 g (「leave」) を生成することができる。そのような AI ベースのトークンジェネレータ 4 1 の一例は、たとえば J. Devlin 他、「BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (BERT: 言語理解のためのディープ双方向トランスフォーマの事前トレーニング)」、arXiv: 1810.04805 に記載されているような、言語の Bidirectional Encoder Representation From Transformers (BERT) モデルを実装する。いくつかの実施形態では、トレーニング中に、図 3 ~ 4 に示されるようにトークンジェネレータ 4 1 がシーケンスアナライザ 4 2 に接続され、ジェネレータ 4 1 が、シーケンスアナライザ 4 2 の調節可能パラメータとは別個の調節可能な機能パラメータ (たとえば、シナプス重みなど) のセットを有し、トレーニング済みジェネレータ 4 1 の少なくとも 1 つの調節可能パラメータが、シーケンスアナライザ 4 2 の出力に従って調整されるという意味で、トークンジェネレータ 4 1 は、シーケンスアナライザ 4 2 と共にトレーニングされる。

#### 【0027】

[0037] 入力修正器 4 0 によって実装される別の例示的シーケンス変換 3 0 は、トークン埋込みベクトルの操作を含み得る。そのような実施形態では、修正後シーケンス 3 4 は、トークンのシーケンス自体ではなく、埋込みベクトルの配列を含み得る。入力修正器 4 0 は、当技術分野で一般に埋込み空間と呼ばれる抽象多次元ベクトル空間内のトレーニングシーケンス 3 2 の各トークンの位置を示す座標のセットを決定するように構成されたトークンエンコーダを含み得る。座標のそれぞれのセットは、それぞれのトークンに関連するトークン埋込みベクトルを集合的に定義する。図 5 は、例示的トークン埋込み空間 5 0 と、トークン 3 5 a ~ b をそれぞれ表すトークン埋込みベクトル 5 5 a ~ b のセットとを示す。

#### 【0028】

[0038] 例示的埋込み空間は軸のセットによって張られ、各軸は、(たとえば、主成分 / 特異値分解実施形態において) 別個のトークン特徴、またはトークン特徴の一次結合を

表す。コンピューティングイベントのシーケンス中の異常を検出するように構成された実施形態では、トークン特徴は、各イベントの様々なイベント特徴（たとえば、イベントのタイプ、経路標識、ネットワークアドレスなど）を含み得る。好ましい実施形態では、トークンがトレーニングシーケンス内のその位置に従って埋め込まれ、言い換えれば、そのコンテキストに従って埋め込まれる。そのようなケースでは、埋込み空間 50 は、抽象コンテキスト空間を含み得、類似のコンテキストで主に生じる 2 つのトークンが、相対的に近くに共に配置される。とりわけ `word2vec`、`GloVe`、および `BERT` を含む、いくつかのそのような埋込みは当技術分野で周知である。埋込みベクトル表現 55 a ~ b を生成するために、トークンエンコーダは、トークンシーケンスのコーパスに関してトレーニングされなければならない、トークンシーケンスのコーパスは、トレーニングコーパス 18、すなわちシーケンスアナライザ 42 をトレーニングするために使用されるコーパスと一致し得る。トレーニングは、当技術分野で周知の任意の方法に従って、たとえば `bag-of-words` および / または `skip-gram` アルゴリズムに従って進行し得る。いくつかの実施形態では、トークンエンコーダの調節可能パラメータがシーケンスアナライザ 42 の出力に従って調整されるという意味で、トークンエンコーダはアナライザ 42 と共にトレーニングされる。

10

#### 【0029】

[0039] 図 6 に示されるように、トレーニングシーケンス 32 をひずませるためのいくつかのシーケンス変換 30 (図 3) は、埋込みベクトルに対して直接的に作用し得る。例示的変換  $T_j$  は、トレーニングシーケンス 32 の選択されたトークンを表す元のトークン埋込みベクトル 55 c を修正後ベクトル 55 d に変更し得る。例示的埋込み変換は、軸のうちの一つに沿って、または変換特有の所定の方向に沿って、小さい量だけベクトルをナッジすることを含む。別の例示的変換は、所定の平面の周りの回転および反射を含み得る。それぞれの変換は、トレーニングシーケンス 32 のすべてのトークンに適用され、あるいはたとえば (上記で示したような) バイナリマスクまたは他の選択基準によって識別される、選択されたトークンだけに適用され得る。

20

#### 【0030】

[0040] いくつかの実施形態では、シーケンスアナライザ 42 は、入力トークンシーケンスを処理して、入力トークンシーケンスに従って決定された変換予測標識 36 およびトークン予測標識 38 を含む予測標識のセットを生成するように構成される。変換予測標識 36 は、入力トークンシーケンスを生成するためにどのシーケンス変換 30 が使用された可能性が高いかを示す。例示的实施形態では、変換予測標識 36 は、複数の数値スコア  $P(T_1)$ 、 $P(T_2)$ 、...  $P(T_k)$  を含み、各スコア  $P(T_j)$  は、それぞれの入力トークンシーケンスを生成するためにそれぞれの変換  $T_j$  が適用された可能性を示す。たとえば、標識 36 は、入力修正器 40 によって実装されるそれぞれの別個のシーケンス変換 30 について別個のスコアを含み得る。スコア  $P(T_j)$  は、所定の境界の間 (たとえば、0 から 1 の間) でスケールされ得、大きい値は高い可能性を示す。

30

#### 【0031】

[0041] 次に、トークン予測標識 38 のいくつかの実施形態は、入力シーケンスのどのトークンが入力修正器 40 によって修正された可能性が高いかを示す。例示的实施形態では、トークン予測標識 38 は複数の数値スコア  $S_1$ 、 $S_2$ 、... を含み、スコア  $S_n$  は、入力シーケンスの  $n$  番目のトークンが入力修正器 40 によって変更された可能性を示す。図 3 に示されるようなトレーニングプロセスでは、各スコア  $S_n$  は、修正後トークンシーケンス 34 の  $n$  番目のトークンがそれぞれのトレーニングシーケンス 32 の  $n$  番目のトークンとは異なる可能性を示し得る。

40

#### 【0032】

[0042] 直感的な視点からは、変換予測標識 36 は、トレーニングシーケンス 32 をひずませるためにどの変換が使用されたかを推測する試みを表し、トークン予測標識 38 は、個々のどのトークンが破壊されたかを推測する試みを表す。標識 36 および 38 は冗長な情報を伝達するように見えるが (結局のところ、各変換は特定のトークンマスクを有す

50

る)、標識 36 および 38 は、シーケンスアナライザ 42 の別個のサブシステム(たとえば、ディープニューラルネットワークのニューロンの別個のグループ)によって生成されるという意味で無関係である。さらに、シーケンス変換 30 とその関連するトークンマスクとの間の結び付きまたは関連の、シーケンスアナライザ 42 内に構築された事前の知識はない。その代わりに、アナライザ 42 は、トレーニング中に自動的にそのような関連を学習し得る。いくつかの実施形態は、標識 36 と 38 の両方を使用すると、シーケンスアナライザ 42 のトレーニングをかなり促進し得、たとえば学習を加速し、標識 36 および 38 のうちのただ 1 つを使用するのと同様の異常検出性能達成するためにかなり小さいトレーニングコーパスの使用を可能にするという観測に依拠する。

#### 【0033】

[0043]シーケンスアナライザ 32 の例示的アーキテクチャが図 7 に示されており、層 / ニューラルネットワークモジュールのスタックを備え、各層は、前の層 / モジュールの出力を受け取り、入力をスタックの次の層に与える。シーケンスアナライザ 32 は、トークン表現 48 の配列の形で入力を受け取り得、各トークン表現は、入力シーケンスのそれぞれのトークンを特徴付ける数のベクトルを含む。トレーニングシナリオでは、図 7 の各トークン表現 48 は、修正後シーケンス 34 (図 3) の別個のトークンを表し得る。当技術分野で 1 ホット符号化と呼ばれるものを用いる例示的一実施形態では、トークン表現 48 は  $N \times 1$  ベクトルを含み、各行は別個のトークンタイプを表し、 $N$  はトークン語彙のサイズを表し、非ゼロ要素は、それぞれのトークンがそれぞれのトークンタイプであることを示す。トークン語彙の例示的サイズ  $N$  は、特定の適用分野では数百から数百万の範囲である。別の例では、各トークン表現 48 は、当技術分野で周知の埋込みアルゴリズムに従って生成されるトークン埋込みベクトルを含み得る。たとえば、上記で論じられた図 5 の埋込みベクトル 55 a ~ b を参照されたい。

#### 【0034】

[0044]それぞれの連続する層  $L_i$  は、それぞれの層に特有のパラメータ(たとえば、活性化、重み、バイアス)のセットに従って、前の層から受け取った入力を変換し、内部ベクトル 49 を生成し、そのサイズおよび値の範囲は、アナライザ 32 の別個の層 / モジュールの間で様々であり得る。たとえば、いくつかの層は、プーリングまたは損失層のケースと同様に、それぞれの入力ベクトルの次元削減を達成する。各層のタイプおよびアーキテクチャは、各実施形態にわたって異なり得る。

#### 【0035】

[0045]シーケンスアナライザ 42 の例示的アーキテクチャは、畳み込みニューラルネットワーク(CNN)層と、その後続く、整流器(たとえば、ReLU または他の活性化関数)および / または損失層にさらに結合された密層(すなわち、全結合層)とを含む。代替実施形態は、再帰型ニューラルネットワーク(RNN)内に供給される CNN 層と、その後続く全結合層および ReLU / 損失層とを含み得る。畳み込み層は、効率的に、内部ベクトル 49 を、当技術分野でフィルタと呼ばれる重みの行列と乗算し、埋込みテンソルを生成し、その結果、それぞれのテンソルの各要素は、選択されたトークンからの寄与を有するが、選択されたトークンに隣接する別のトークンからの寄与も有する。したがって、埋込みテンソルは、個々のトークンよりも粗いグラニューラリティで、入力トークンシーケンスを集合的に表す。フィルタ重みは、トレーニングプロセス中に調節され得る調節可能パラメータである。

#### 【0036】

[0046]再帰型ニューラルネットワーク(RNN)は、人工ニューラルネットワークの特別なクラスを形成し、ネットワークノード間の結び付きが有向グラフを形成する。とりわけ長・短期記憶(LSTM)ネットワークおよびグラフニューラルネットワーク(GNN)を含む RNN のいくつかのフレーバが当技術分野で周知である。典型的な RNN は隠れユニット(たとえば、個々のニューロン)のセットを含み、ネットワークのトポロジは、各隠れユニットがそれぞれのトークン  $m_j$  を特徴付ける入力(たとえば、埋込みベクトル)だけでなく、隣接する隠れユニットによって与えられる入力も受け取るように具体的

10

20

30

40

50

に構成され、隣接する隠れユニットは、入力トークンシーケンス内のトークン  $m_j$  に先行するトークン  $m_{j-1}$  を特徴付ける入力を受け取る。その結果、各隠れユニットの出力は、それぞれのトークン  $m_j$  によって影響を受けるだけでなく、先行するトークン  $m_{j-1}$  によっても影響を受ける。言い換えれば、RNN層は、前のトークンのコンテキスト内の各トークンについての情報を処理し得る。双方向RNNアーキテクチャは、入力トークンシーケンスの前のトークンと後続のトークンの両方のコンテキスト内の各トークンについての情報を処理し得る。

【0037】

[0047]シーケンスアナライザ42のさらに別の例示的实施形態は、トランスフォーマニューラルネットワーク層のスタックを備え得る。トランスフォーマアーキテクチャは、とりわけ、たとえばA. Vaswani等「Attention is all you need (あなたが必要なすべてのものは注意である)」、arXiv:1706.03762に記載されている。各入力トークンシーケンスについて、トランスフォーマ層は、コンテキスト化されたトークン埋込みベクトルのシーケンスを生成し得、各トークン埋込みベクトル  $h_j$  は、入力シーケンスの複数の(たとえば、すべての)トークン  $m_j$  からの情報を符号化する。トランスフォーマ層の出力は、当技術分野で予測ヘッドと呼ばれ、図7のブロックH1およびH2として示される複数の別個の分類器モジュール(たとえば、密層)内に供給され得る。次いで、ヘッドH1およびH2は、それぞれ変換予測標識36およびトークン予測標識38を出力し得る。

10

【0038】

[0048]図8は、本発明のいくつかの実施形態による、異常検出器20をトレーニングするためにAIトレーニングアプライアンス14(図1)によって実施されるステップの例示的シーケンスを示す。コーパス18からトレーニングシーケンス32を選択することに対応して、ステップ204~206のシーケンスは、シーケンス変換30の利用可能なセットから変換を(たとえば、ランダムに)選択し、それぞれの変換をトレーニングシーケンス32に適用し、したがって修正後トークンシーケンス34を生成し得る。

20

【0039】

[0049]修正後シーケンス34はシーケンスアナライザ42に供給され、シーケンスアナライザ42は、修正後シーケンス34を処理して、予測標識36および/または38を生成する。次いで、ステップ208は、トレーニングシーケンス32および/または予測標識36および/または38に従ってユーティリティ関数を決定し得る。機械学習の技術分野で損失とも呼ばれる例示的ユーティリティ関数は、以下のように表現され得る。

30

【0040】

$$U(\theta) \quad [1]$$

上式で  $U$  はトレーニングシーケンスを表し、 $\theta$  は調節可能パラメータのセットを表す。トレーニングは、ユーティリティ  $U$  を最小にする方向にパラメータ  $\theta$  を調節することを含み得る。

【0041】

[0050]単純なユーティリティ関数は、所望の出力からのアナライザ42の出力の逸脱を定量化し得る。たとえば、例示的ユーティリティ関数は、ステップ206でどの変換が適用されたか、および/またはステップ206で元のトレーニングシーケンス32のどのトークンが破壊されたかをアナライザ42が正しく推測したかどうかを示し得、誤った推測に対してアナライザ42にペナルティを課し得る。

40

【0042】

[0051]いくつかの実施形態は、標識36および38に従って決定されたユーティリティを組み合わせると、トレーニングを促進し、かつ/またはより高性能のシーケンスアナライザ42が得られ得るといふ観測に依拠する。好ましい実施形態は、シーケンスレベル成分(シーケンス変換の選択を示す)をトークンレベル成分(それぞれの個々のトークンが破壊されたか否かを示す)と組み合わせるアグリゲートユーティリティ関数を使用し得る。

50

【 0 0 4 3 】

$$U = \alpha_1 U_S + \alpha_2 U_T \quad [2]$$

上式で  $U_S$  および  $U_T$  は、それぞれシーケンスレベル成分およびトークンレベル成分を表し、 $\alpha_1$  および  $\alpha_2$  は、各ユーティリティ関数の相対的寄与を変更することを可能にする重みである。最尤トレーニング方法を実装するいくつかの実施形態では、

【 0 0 4 4 】

【数 1】

$$U_S = \mathbb{E}_{k,x} [-\log P(T_k | \tilde{x}, \theta_A)], \quad [3]$$

10

【 0 0 4 5 】

上式で

【 0 0 4 6 】

【数 2】

$$\mathbb{E}$$

20

【 0 0 4 7 】

は期待値を表し、

【 0 0 4 8 】

【数 3】

$$P(T_k | \tilde{x}, \theta_A)$$

30

【 0 0 4 9 】

は、ひずんだシーケンス

【 0 0 5 0 】

【数 4】

$$\tilde{x}$$

40

【 0 0 5 1 】

がシーケンス変換  $T$  の適用によって生成された確率を表し（たとえば、図 3 の変換予測標識 3 6 を参照）、 $\theta_A$  は一般に、シーケンスアナライザ 4 2 の調節可能パラメータを表す。一方、

【 0 0 5 2 】

【数 5】

50



$$U_T = \mathbb{E}_{k,x} [\sum_i -\log S_i (T_k | \tilde{x}, \theta_A)], \quad [4]$$

【0053】

上式で

【0054】

【数6】

10

$$S_i(T_k | \tilde{x}, \theta_A)$$

【0055】

は、トレーニングシーケンスのトークン*i*がシーケンス変換Tの適用によって影響を受けた確率を表す（たとえば、図3の変換予測標識38を参照）。

[0052]いくつかの実施形態では、入力修正器40の構成要素（トークン埋込みベクトルを生成するように構成されたトークンジェネレータ41および/またはトークンエンコーダなど）が、シーケンスアナライザ42と共にトレーニングされる。そのような実施形態は、前述のU<sub>S</sub>およびU<sub>T</sub>に加えて、ジェネレータユーティリティ関数を使用し得る。

20

【0056】

$$U = \alpha_1 U_S + \alpha_2 U_T + \alpha_3 U_G \quad [5]$$

上式で、 $\alpha_3$ は、ジェネレータユーティリティ関数U<sub>G</sub>のグローバルユーティリティへの寄与を調節するために使用される別の重みを表し、

【0057】

【数7】

30

$$U_G = \mathbb{E}_{k,x} [\sum_i -\log P_G (t_i | \tilde{x}, \theta_G)], \quad [6]$$

【0058】

上式で  $\theta_G$  は一般に、トークンジェネレータ41の調節可能パラメータを表し、

【0059】

【数8】

40

$$P_G (t_i | \tilde{x}, \theta_G)$$

【0060】

は、トークン  $t_i$  が修正後シーケンス34中に現れる確率、または言い換えれば、トークン  $t_i$  が修正後シーケンス34のコンテキスト内で妥当である確率を表す。

[0053]次いで別のステップ210は、決定されたユーティリティ関数に従って、パラメータ  $\theta_A$  および/または  $\theta_G$  のセットを調節し得る。そのような調節は、勾配降下によ

50

る逆伝播手順、または選ばれたユーティリティ関数を最小にすることを目標とする任意の他の最尤探索を実装し得る。トレーニングは、終了条件が満たされるまで、たとえば所定のエポック数について、所定の数のトレーニングシーケンスが解析されるまで、所定のレベルの異常検出性能が実証されるまでなどの間、続行し得る（ステップ 212）。トレーニングの成功に回答して、シーケンスアナライザ 42 の調節可能パラメータ（たとえば、シナプス重みなど）の最適な値が、検出器パラメータ値 24（図 2）の形でエクスポートされ、クライアントシステム 10 a ~ c に送られる。

【0061】

[0054] 図 9 は、本発明のいくつかの実施形態による、異常を検出するためにクライアントシステム 10 a ~ c および / またはユーティリティサーバ 12 によって実施されるステップの例示的シーケンスを示す。ステップ 222 ~ 224 のシーケンスは、AI トレーニングアプライアンス 14（図 2）から検出器パラメータ値 24 を受け取り、それぞれの値で検出器 20 のローカルインスタンスをインスタンス化することによって、動作のために検出器 20 を準備する。次いで、各ターゲット・トークン・シーケンス 22 について、ステップ 228 は、シーケンスアナライザ 42 を実行して、それぞれのターゲット・トークン・シーケンス 22 についてのトークン予測標識 38 を決定し得る。

10

【0062】

[0055] 別のステップ 230 は、判断モジュール 44 を適用して異常標識 26 を生成し得る。いくつかの実施形態では、判断モジュール 44 は、たとえばターゲット・トークン・シーケンス 22 全体にわたって取られた個々のトークン予測スコア  $S_i$  の平均として、トークン予測標識 38 に従って異常標識 26 を決定するように構成される。それぞれの個々のスコア  $S_i$  は、図 3 に関連して上記で説明されたように、シーケンス 22 のそれぞれのトークンがシーケンス変換 30 の適用によって破壊されたかどうかの可能性を定量化し得る。いくつかの実施形態は、以下に従ってシーケンス特有の異常スコア  $A$  を決定し得る。

20

【0063】

【数 9】

$$A = \frac{1}{L_S} \sum_i S_i, \quad [7]$$

30

【0064】

[0056] 上式で、 $L_S$  は、ターゲットシーケンス 22 の長さ（トークンのカウント）を表す。大きい  $S_i$  が、それぞれのトークンが破壊された可能性が高いことを示す一実施形態では、 $A$  の大きい値は、ターゲット・トークン・シーケンス 22 が異常である可能性が高いことを示し得る。逆に、大きい  $S_i$  が、それぞれのトークンが破壊されていない可能性が高いことを示すとき、大きい  $A$  の値は、ターゲットシーケンス 22 が異常ではないことを示し得る。判断モジュール 44 のいくつかの実施形態は、異常スコア  $A$  の計算された値を所定のしきい値と比較し、比較の結果に従ってターゲット・トークン・シーケンス 22 が異常であるか否かを判定する。

40

【0065】

[0057] 図 10 は、本明細書で説明される方法のうちのいくつかを実行するようにプログラムされたコンピューティングアプライアンス 70 の例示的ハードウェア構成を示す。コンピューティングアプライアンス 70 は、図 1 のクライアントシステム 10 a ~ c、ユーティリティサーバ 12、および AI トレーニングアプライアンス 14 のいずれかを表し得る。図示されるコンピューティングアプライアンスはパーソナルコンピュータであり、

50

サーバ、携帯電話、タブレットコンピュータ、ウェアラブルなどの他のデバイスは、わずかに異なる構成を有し得る。プロセッサ72は、信号および/またはデータのセットと共に計算および/または論理演算を実行するように構成された物理デバイス(たとえばマイクロプロセッサ、半導体基板上に形成されたマルチコア集積回路)を備え得る。そのような信号またはデータは、プロセッサ命令、たとえば機械コードの形で符号化され、プロセッサ72に送達され得る。

#### 【0066】

[0058]プロセッサ72は一般に命令セットアーキテクチャ(ISA)によって特徴付けられ、ISAは、とりわけ、プロセッサ命令のそれぞれのセット(たとえばx86ファミリとARM(登録商標)ファミリ)、およびレジスタのサイズ(たとえば、32ビットプロセッサと64ビットプロセッサ)を指定する。プロセッサ72のアーキテクチャは、所期の主な用途に従って様々であり得る。中央演算処理装置(CPU)は汎用プロセッサであり、グラフィックス処理装置(GPU)は、イメージ/ビデオ処理およびいくつかの形態の並列コンピューティングのために最適化される。プロセッサ72は、Google(登録商標), Inc.によるテンソル処理装置(TPU)、様々な製造業者によるニューラル処理装置(NPU)などの特定用途向け集積回路(ASIC)をさらに含み得る。TPUおよびNPUは、本明細書で説明されるような機械学習アプリケーションに特に適していることがある。

#### 【0067】

[0059]メモリユニット74は、動作を実施する過程でプロセッサ72によってアクセスされ、または生成されるデータ/信号/命令符号化を記憶する揮発性コンピュータ可読媒体(たとえば、ダイナミックランダムアクセスメモリ-DRAM)を備え得る。入力デバイス76は、とりわけ、ユーザがデータおよび/または命令をアプライアンス70内に導入することを可能にするそれぞれのハードウェアインターフェースおよび/またはアダプタを含む、コンピュータキーボード、マウス、およびマイクロフォンを含み得る。出力デバイス78は、とりわけモニターやスピーカなどのディスプレイデバイス、ならびにそれぞれのコンピューティングアプライアンスがユーザにデータを通信することを可能にする、グラフィックカードなどのハードウェアインターフェース/アダプタを含み得る。いくつかの実施形態では、入力および出力デバイス76~78は、共通のハードウェア(たとえば、タッチスクリーン)を共有する。記憶デバイス82は、ソフトウェア命令および/またはデータの非揮発性記憶、読取り、および書込みを可能にするコンピュータ可読媒体を含む。例示的記憶デバイスは、磁気ディスク、光ディスク、およびフラッシュメモリデバイス、ならびにCDおよび/またはDVDディスクおよびドライブなどの取外し可能媒体を含む。ネットワークアダプタ84は、コンピューティングアプライアンス70が電子通信ネットワーク(たとえば、図1のネットワーク15)および/または他のデバイス/コンピュータシステムに接続することを可能にする。

#### 【0068】

[0060]コントローラハブ80は一般に、複数のシステム、周辺機器、および/またはチップセットバス、ならびに/あるいはプロセッサ72とアプライアンス70のハードウェア構成要素の残りの部分との間の通信を可能にするすべての他の回路を表す。たとえば、コントローラハブ80は、メモリコントローラ、入力/出力(I/O)コントローラ、および割込みコントローラを備え得る。ハードウェア製造業者に依りて、いくつかのそのようなコントローラは単一の集積回路内に組み込まれ得、かつ/またはプロセッサ72と一体化され得る。別の例では、コントローラハブ80は、プロセッサ72をメモリ74に接続するノースブリッジ、ならびに/あるいはプロセッサ72をデバイス76、78、82、および84に接続するサウスブリッジを備え得る。

#### 【0069】

[0061]前述の例示的システムおよび方法は、様々な適用分野で異常の効率的な自動検出を可能にする。いくつかの実施形態では、トレーニングコーパスから引き出されるトークンシーケンスが、トレーニングを施すシーケンスアナライザに供給される前に、複数の

10

20

30

40

50

所定のシーケンス変換のうちの少なくとも1つに従ってひずめられる。次いでシーケンスアナライザは、それぞれの入力トークンシーケンスを生成するためにどの変換が使用されたかを正しく推測するようにトレーニングされる。

【0070】

[0062]異常検出器をトレーニングするためのいくつかの従来の手順は、トークンのうちのいくつかをランダムに置き換えることによってトレーニング・トークン・シーケンスを破壊し、その後で、どのトークンが置き換えられたかを推測するように検出器をトレーニングする。しかしながら、そのようなトレーニング方法は、計算資源の点で比較的成本がかかり、トレーニングコーパスおよび/またはアプリケーションのいくつかの選択肢については不安定であり得る。この従来手法とは対照的に、本発明のいくつかの実施形態は、変換の事前定義されたセットを使用して入力トークンシーケンスをひずませ、トークンレベル成分（それぞれの個々のトークンが破壊されたか否かを示す）をシーケンスレベル成分（入力シーケンス全体をひずませる方式を示す）と組み合わせるアグリゲートユーティリティ関数に従ってトレーニングする。より従来型のトークンレベルユーティリティに加えてシーケンスレベルユーティリティを使用することは、反直感的に見えることがある。いくつかの実施形態では、各シーケンス変換が特定のトークンマスクを有し、したがってどの変換が適用されたかを推測することが、実質的にはどのトークンが破壊されたかも推測することになり得るからである。しかしながら、いくつかの実施形態は、シーケンスレベルの学習タスク（適用された変換を推測すること）をセットアップすると同時に、トークンレベルのタスク（特定のトークンが破壊されたかどうかを推測すること）をセットアップすることは、シーケンスアナライザが存在しないマスクパターンを予測するのを阻止することによって、正しい学習を強化し得る。いくつかの実施形態では、トークンレベルの予測およびシーケンスレベルの予測が、ディープニューラルネットワークの別個の予測ヘッドによって生成される。したがって、異常検出器内に構築される変換とトークンマスクとの間の相関の事前の知識はない。その代わりに、検出器は、トレーニング中にそのような潜在的相関を学習し、そのことは、よりロバストなモデルに至り得る。

10

20

【0071】

[0063]コンピュータ実験は、トークンレベルのタスクとシーケンスレベルのタスクを組み合わせると、検出器の性能を改善することによって学習を促進することを示している。逆に、本明細書で説明されるようなトレーニング方法を使用することにより、より小さいトレーニングコーパスおよび/またはより少ないネットワークパラメータを使用して、同一のレベルの異常検出性能が達成され得る。これは、トレーニングコーパス18が比較的小さいサイズを有する状況（たとえば、トレーニングコーパスがソーシャルメディアポストからなるとき）での作成者属性などの異常検出タスクにとって特に有利であり得る。本明細書で説明されるように検出器をトレーニングすることは、直感的にはトレーニングコーパスのサイズを人工的に増加させることに対応する。同一のトレーニング・トークン・シーケンス32が、別個のシーケンス変換30の適用に回答して複数の別個の修正後シーケンス34を引き起こし得るからである（図3参照）。

30

【0072】

[0064]いくつかの実施形態は、第2のAIシステムを利用して、トレーニングシーケンスの妥当なひずみを生成する。たとえば、それぞれのトークンシーケンスの残りの部分のコンテキストが与えられると、BERT言語モデルを実装するトークンジェネレータが使用され、選択されたトークンが妥当な置換で置き換えられ得る。いくつかの実施形態は、既にトレーニングされたジェネレータが、異常と見なされるにはある意味で「妥当過ぎる」修正後トレーニングシーケンスを生成することによって学習を妨げ得るという観測に依拠して、事前トレーニング済みの高性能バージョンのトークンジェネレータを使用する代わりに、トークンジェネレータを異常検出器と共に明示的にトレーニングする。共にトレーニングすることは、シーケンスアナライザが修正を検出することにより熟練するにつれて、トークンジェネレータが妥当な修正後トレーニングシーケンスを生成することにますます熟練することを保証し得る。さらに、異常検出器をトレーニングするために使用さ

40

50

れるコーパスとは別個のコーパスに関してトークンジェネレータを事前トレーニングすることは、異常値情報をもたらし、異常検出器が異常値情報をそのように認識することを防止し得る。

#### 【0073】

[0065]本発明のいくつかの実施形態に従ってトレーニングされた異常検出器は、とりわけ、以下を含む様々なシナリオで使用され得る。

#### 自動テキスト分類

[0066]例示的自然言語処理(NLP)アプリケーションでは、異常検出器20は、特定のカテゴリ(たとえば、ビジネスニュース)に属するテキストからなるコーパスに関してトレーニングされ、次いでターゲット・テキスト・フラグメントがそれぞれのカテゴリに属するか否かを判定するために使用され得る。そのような実施形態では、高い異常スコアは、それぞれのテキストがそれぞれのカテゴリに属さないことを示し得る。

10

#### 【0074】

[0067]コンピュータ実験では、本明細書で説明される異常検出器が、ニュース記事の標準基準コーパス(20 News groups)のサブセットに関してトレーニングされ、サブセットは、選択されたカテゴリ(コンピューティング、リクリエーション、科学、雑報、政治、または宗教)の記事からなる。実験は、4つの積重ねトランスフォーマ層を含み、2つの予測ヘッドが上端にあるシーケンスアナライザを使用した。各トランスフォーマ層は、4つの自己注意ヘッド、サイズ256の隠れ層、ならびにサイズ1024および256のフィードフォワード層を含んだ。各予測ヘッドは、非線形性によって分離される2つの線形層を有し、分類層で終了した。トレーニング・トークン・シーケンスの最大サイズは128であった。シーケンス変換30は、ランダムトークンジェネレータを使用して、別個のマスクパターンに従ってトークンを置き換えることからなるものであった。様々なカウントおよびカバレッジのマスクパターンが、入力トレーニングシーケンスの25%から50%の間をカバーする、5から100個の間の別個のマスクパターンと共に試行された。

20

#### 【0075】

[0068]次いで、トレーニングされた検出器は、集合からランダムに選択された記事が、検出器がトレーニングされたカテゴリに属するかどうかを識別するように求められた。本発明のいくつかの実施形態に従ってトレーニングされた異常検出器は、それぞれのタスクで現況技術の従来の異常検出器より一貫して大幅に優れており、典型的なarea under the receiver operating curve(AUROC)値は、約70%(科学カテゴリに関してトレーニングされたとき)から92%超(コンピューティングニュースに関してトレーニングされたとき)までの範囲である。一般には、別個の変換数を増加させると、トレーニング済みの異常検出器の性能を改善し、トークン埋込み内の表現性を促進することを実験は明らかにした。25%から50%の割合の破壊されたトークンを有する変換が最良の結果を生み出すように見えた。

30

#### 自動作成者属性

[0069]異常検出器20のいくつかの実施形態は、選択された作成者によって書かれたテキスト(たとえば、手紙、記事、ブログポスト、eメール、ソーシャルメディアポスト)のコーパスに関してトレーニングされ、次いで、ターゲット・テキスト・フラグメントがそれぞれの人によって作成されたかどうかを判定するために使用され得る。例示的アプリケーションは、匿名の手紙の作成者を判定し、様々な文書の真正性、および文学作品の死後の属性を検証することを含む。いくつかの実施形態はまた、犯罪学の適用分野をも有し得る。たとえば、Dark Webリソースの作成者またはユーザを識別する、たとえば盗品、クレジットカードデータ、小児ポルノ、銃、麻薬などの取引などの犯罪活動に関連するユーザについての集合地点として役目を果たすフォーラム上でポストされたメッセージの作成者を識別する際に法執行に関心が持たれることがある。「Dark Web」という用語は、本明細書では、検索エンジンによって索引付けされず、かつ/またはプライベートピアツーピアネットワークもしくはTorなどの匿名化ソフトウェアを介しての

40

50

みアクセス可能であるコンテンツを表すために使用される。

【0076】

[0070]次いで、容疑者のセットによって作成され、公に利用可能であるオンラインコンテンツ（たとえば、人気のあるソーシャルメディアサイトおよび/またはユーザフォーラム上でそれぞれの容疑者によってポストされたコメント）のコーパスに関してトレーニングされた異常検出器の例示的实施形態が、Dark Webから取り入れられたターゲット・テキスト・フラグメントを解析するために使用され得る。ターゲットテキストが異常ではないことを示す異常スコアが、ターゲットテキストの作成者が、検出器がトレーニングされたテキストのコーパスの作成者のうちの1人と一致することを示し得る。

ソーシャルメディア監視

[0071]異常検出器20の一実施形態が、ソーシャルメディアアカウントの選択されたセットに関連するウェブコンテンツ、たとえばTwitter（登録商標）フィードの特定の集合に関してトレーニングされ得る。トレーニングコーパスは、特定の時間ウィンドウ（たとえば、1日、1週など）以内に発行されたコンテンツにさらに限定され得る。次いで、検出器が使用され、新しくポストされたコンテンツが解析され得る。異常は、トピックの変化および/または進行中の交換のトーンの変化を示し得、したがって新しいトピックおよびトレンドの適時の自動検出を可能にする。

フェイクおよび自動生成されたコンテンツの検出

[0072]異常検出器20の一実施形態は、選択された人間の作成者（たとえば、実際のニュース記事、実際のユーザによってポストされたソーシャルメディア）によって書かれたテキストのコーパスに関してトレーニングされ得る。コーパスは、選択された定期刊行物、新聞、またはニュースウェブサイトについて書かれた記事、または選択されたジャーナリストによって書かれた記事にさらに狭められ得る。次いで、トレーニングされた異常検出器は、ターゲット・テキスト・フラグメントを解析するために使用され得る。ターゲットテキストが異常であることを示す異常スコアは、それぞれのテキストがフェイクニュースを含み得ること、および/またはマシンで生成され得ることを示し得る。

データ保護およびプライバシー

[0073]いくつかのクラウドコンピューティングサービスは、ユーザが他のユーザと共有するために、または様々な操作（たとえば、マルウェア走査）のためにリモートサーバにファイルをアップロードすることを可能にする。一例として、ユーザのコンピュータ上で実行中のソフトウェアエージェントは、それぞれのユーザによって示され得る、選択されたフォルダのコンテンツを自動的にアップロードし得る。クラウドにデータをアップロードすることは、ユーザがアップロード用のコンテンツを明示的に選ばないときは特に、プライバシーリスクを含み得る。たとえば、ユーザが何らかの機密データ（たとえば、個人ファイルまたは写真、医療記録など）をアップロードフォルダに誤ってドロップした場合、それぞれのデータがユーザの希望に反して自動的にアップロードされる。

【0077】

[0074]異常検出器20の一実施形態が、ユーザのコンピュータ上にインストールされ、通常はそれぞれのユーザによってアップロードされるファイル、たとえばリモート走査のために最近アップロードされた100個のファイルに関してトレーニングされ得る。追加のフィルタが、Portable Document Format（PDF）文書やMicrosoft（登録商標）Office（登録商標）ファイルなどの特定の種類のファイルのみを選択し得る。そのような実施形態は、前述のような自然言語処理技法を使用し得、トークンは個々の語などを含む。次いで、トレーニング済みの異常検出器が、アップロードに備えてリモート走査のために現在目印が付けられている各ファイルを解析するために使用され得る。それぞれのファイルについて決定された異常スコアが、潜在的異常を示すとき、いくつかの実施形態は、それぞれのファイルがアップロードされるのを防止し得、ユーザに通知し得る。

コンピュータセキュリティ

[0075]異常検出器20のいくつかの実施形態は、活動の正常パターンを表すと見なさ

10

20

30

40

50

れる基準時間間隔中に生じるコンピューティングイベントのシーケンスに関してトレーニングされ、次いでセットクライアントコンピュータシステムの挙動を監視するために使用され得る。クライアント上で検出された異常な挙動は、コンピュータセキュリティ脅威、たとえばそれぞれのクライアントが悪意のあるソフトウェアを実行していること、または侵入者/ハッカーがそれぞれのクライアントへのアクセス権を得たことを示し得る。

#### 【0078】

[0076]いくつかの実施形態では、異常検出は、監視されるソフトウェアエンティティ（たとえば、プロセス、仮想マシンなど）の実行中に生じるイベントのシーケンスを解析することを含む。そのようなイベントの例は、とりわけ、プロセス/スレッドの起動（たとえば、ユーザがアプリケーションを起動する、親プロセスが子プロセスを生み出すなど）、それぞれのクライアントシステムの入力デバイス（たとえば、カメラ、マイクロフォン）にアクセスする試み、ローカルまたはリモートネットワークリソースにアクセスする試み（たとえば、特定のURLのアクセスするためのハイパーテキスト転送プロトコル - HTTP要求、ローカルネットワークを介して文書リポジトリにアクセスする試み）、特定のユニフォームリソース識別子方式で構築された要求（たとえば、mailto:またはFTP:要求）、特定のプロセス命令（たとえば、システムコール）の実行、ライブラリ（たとえば、ダイナミックリンクライブラリ - DLL）をロードする試み、新しいディスクファイルを作成する試み、ディスク上の特定の場所から読み取り、または特定の場所へ書き込む試み（たとえば、既存のファイルを上書きする試み、特定のフォルダまたは文書を開く試み）、および電子メッセージ（たとえば、eメール、ショートメッセージサービス - SMSなど）を送る試みを含む。いくつかの実施形態では、不活動の期間、イベント間の時間ギャップおよび/またはそれぞれのクライアントシステムが遊休状態であり、ユーザ活動を登録せず、または内部システムタスクのみを実施するときの時間間隔も、イベントとしての資格が与えられ得る。本明細書で説明されるシステムおよび方法が、とりわけソーシャルメディア上のユーザの活動に関するイベント、ユーザのブラウジング履歴、ユーザのゲーミング活動などの他の種類のイベントを解析するように適合され得ることを当業者は理解されよう。

#### 【0079】

[0077]イベント検出は、当技術分野で周知の任意の方法を含み得る。一例として、保護されたクライアント上で実行中のセキュリティエージェントは、監視されるソフトウェアエンティティのセットを、event tracking for Windows（登録商標）などのOS 40のイベントロギングサービスに登録し得る。それに応答して、エージェントは、それぞれのプロセスの実行中に生じる様々なイベントの通知を、リアルタイムで、またはログ形式で受信し得る。イベントロギングツールは通常、各イベントについてのタイムスタンプ、イベントタイプを識別する数値コード、それぞれのイベントを生成したプロセスまたはアプリケーションのタイプの標識、および他のイベントパラメータを含むイベント記述子のリストを生成する。イベントシーケンスは、ログを構文解析することによって組み立てられ得る。

#### 【0080】

[0078]いくつかの実施形態は、各イベントを別々のトークンとして扱い得る。トークンはイベント語彙に従って符号化され得、イベント語彙は、数千から数百万個の別個のイベントタイプを含み得る。次いで、異常検出器をトレーニングすることは、前述のようにトレーニングイベントシーケンスに様々な変換を適用することを含み得る。例示的シーケンス変換は、トレーニングシーケンスの選択されたイベントを削除、挿入、および並べ替えること、ならびに選択されたイベントを異なる種類の代替イベントと置き換えることを含み得る。

#### 【0081】

[0079]代替実施形態は、イベントログエントリをテキストトークンのシーケンスと見なし得る。たとえば、ログエントリ：

20:10 | INFO | manager.strage | Found block r

10

20

30

40

50

dd\_\_2\_\_3 locally

が以下のトークンシーケンスに構文解析され得る。

【0082】

20:10;INFO;manager;storage;Found;block;  
rdd\_\_2\_\_3;locally

ただし個々のトークンがセミコロンで分離される。次に、入力修正器40が、選択されたトークンを代替と置き換えることによって、それぞれのトークンシーケンスをひずませ得る。置換のために選択されたトークンの位置が、前述のようにマスクによって示され得る。そのような一例として、上記について決定された修正後トークンシーケンスは以下のように理解され得る。

【0083】

20:10;DEBUG;manager;thread;Found;block;  
rdd\_\_2\_\_3;globally

ただし、代替トークンが太字で示されている。いくつかの実施形態では、プログエントリの選択されたフィールドが修正されないようにマスクが選ばれる。トークンジェネレータ41は、候補のフィールド特有のプールまたは位置特有のプールから代替トークンを選択するように構成され得る。上記の例では、第2のトークンについての候補代替のプールが{WARNING, DEBUG, INFO, CRITICAL}からなることがある。

【0084】

[0080]ログに対する異常検出の例示的適用は、ハニーポットシステム上に記録されたアクセスおよび/またはイベントログを解析することによってゼロデイ攻撃を検出することを含む。本明細書で説明されるような異常検出器は、ログの第1の部分に関してトレーニングされ得、したがってログの第1の部分に対応する時間枠の間のそれぞれのハニーポットの「正常な」挙動を学習する。次いで、異常検出器は、ログの第2の部分解析するために使用され得る。異常がログの第1の部分と第2の部分との間のハニーポットの挙動の変化を示し得、新しいマルウェアの可能な出現、ポットネットの活動化などを示唆する。いくつかの実施形態は、異常検出器を周期的に(たとえば、前の時間からのログデータに関して毎時間ごとに)再トレーニングし、それを使用して新しい脅威をリアルタイムに監視し得る。

【0085】

[0081]本発明の範囲から逸脱することなく、上記の実施形態が多くの方で変更され得ることは、当業者にとっては明らかであろう。したがって、本発明の範囲は、以下の特許請求の範囲およびその法的均等物によって決定されるべきである。

【図面】

【図1】

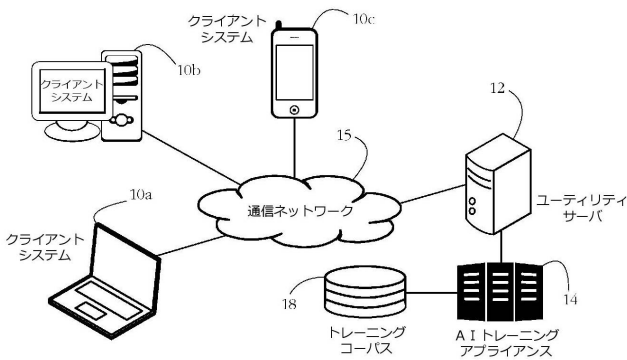


FIG. 1

【図2】

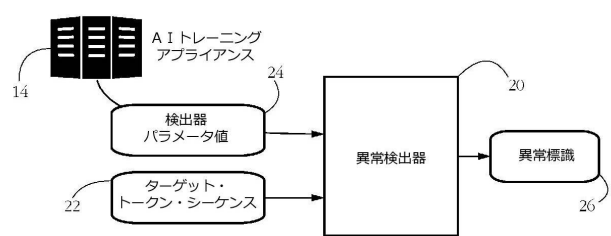


FIG. 2

10

20

30

40

50



【 図 3 】

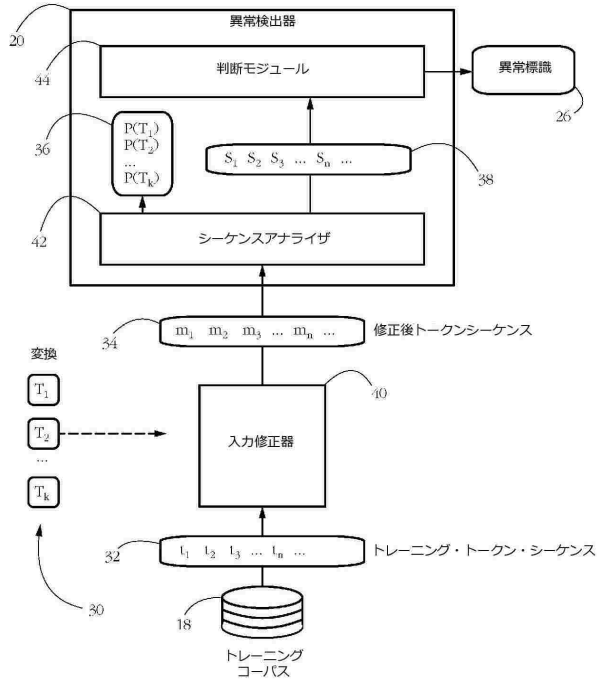


FIG. 3

【 図 4 】

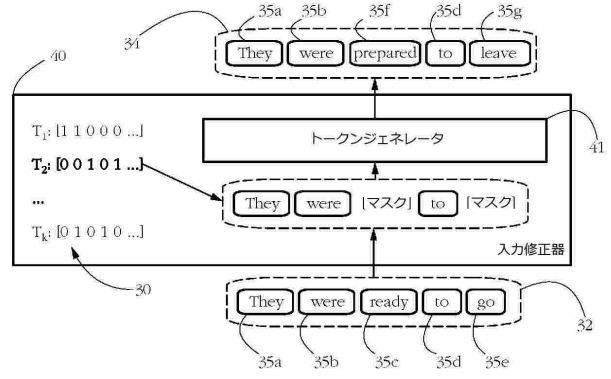


FIG. 4

【 図 5 】

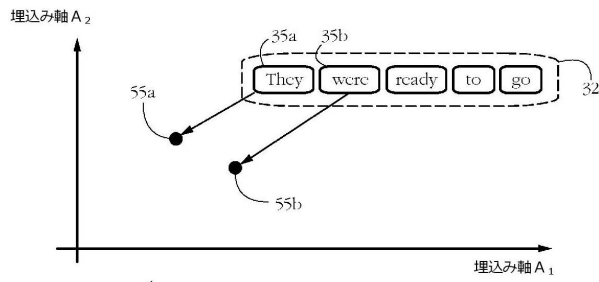


FIG. 5

【 図 6 】

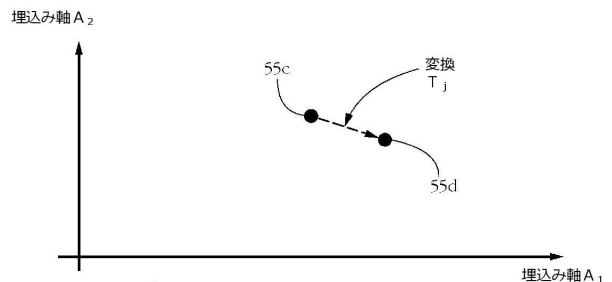


FIG. 6

10

20

30

40

50

【 図 7 】

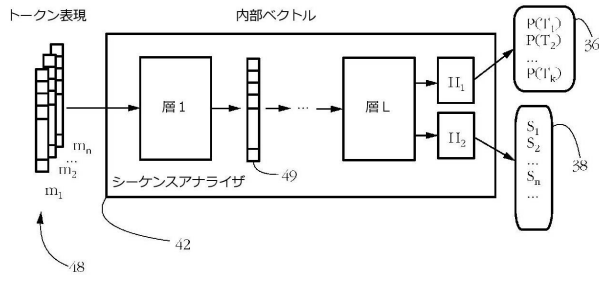


FIG. 7

【 図 8 】

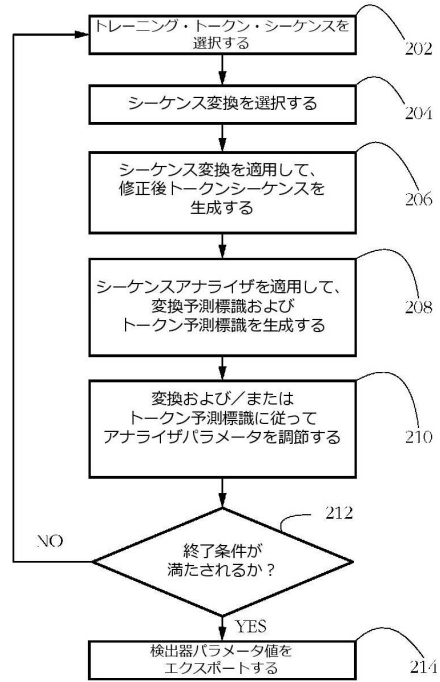


FIG. 8

【 図 9 】

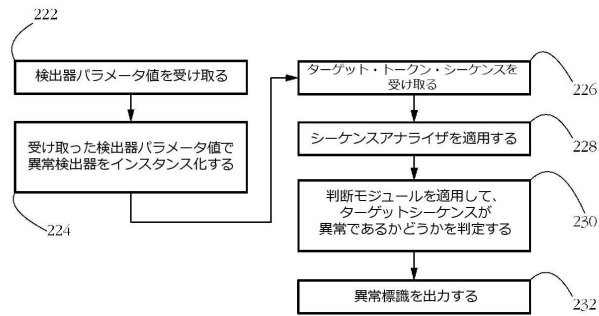


FIG. 9

【 図 10 】

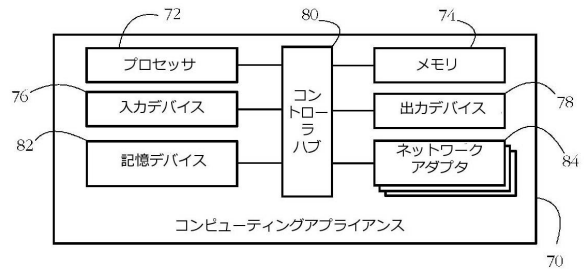


FIG. 10

10

20

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No  
PCT/EP2022/058130

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> INV. G06N3/04 G06N3/08 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) <b>G06N</b>		
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) <b>EPO-Internal, WPI Data</b>		
<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	<b>WO 2020/255137 A1 (YISSUM RES DEV CO OF HEBREW UNIV JERUSALEM LTD [IL])</b> <b>24 December 2020 (2020-12-24)</b> <b>paragraph [0026] - paragraph [0038]</b> <b>paragraph [0057] - paragraph [0073]</b> <b>figure 1</b>	1-21
A	----- <b>EP 3 726 409 A2 (CROWDSTRIKE INC [US])</b> <b>21 October 2020 (2020-10-21)</b> <b>paragraph [0029] - paragraph [0046]</b> <b>paragraph [0166] - paragraph [0176]</b> -----	1-21
<input 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 <b>23 August 2022</b>		Date of mailing of the international search report <b>08/09/2022</b>
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer <b>Baldan, Marco</b>

Form PCT/ISA/210 (second sheet) (April 2005)

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No

**PCT/EP2022/058130**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
<b>WO 2020255137 A1</b>	<b>24-12-2020</b>	<b>EP 3987455 A1</b>	<b>27-04-2022</b>
		<b>US 2022253699 A1</b>	<b>11-08-2022</b>
		<b>WO 2020255137 A1</b>	<b>24-12-2020</b>
-----			
<b>EP 3726409 A2</b>	<b>21-10-2020</b>	<b>EP 3726409 A2</b>	<b>21-10-2020</b>
		<b>US 2020327225 A1</b>	<b>15-10-2020</b>
-----			

10

20

30

40

50

## フロントページの続き

MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,SM,TR),OA(BF,BJ,CF,CG,CI,CM,GA,GN,GQ,GW,KM,ML,MR,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,IT,JM,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,ST,SV,SY,TH,TJ, TM,TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

トラダ・ミトロポリト・ベニアミン・コスタチェ・ヌマルル 1 ,

(72)発明者

フローリン・エム , ブラッド

ルーマニア国 1 0 5 6 0 0 クンピナ , ジュデトゥル・プラホバ , ストラダ・ペトレ・リシウ・ヌマルル 4

(72)発明者

アレクサンドゥル , ノバク

ルーマニア国 0 4 1 3 3 4 ブカレスト , エスオーエス・オルテニティ・ヌマルル 2 5 4 , ブロック 1 5 1 エスセ 1 , アパルタメントゥル 1 8 , セクター 4

(72)発明者

エレナ , プルチャヌ

ルーマニア国 0 4 0 3 6 5 ブカレスト , ストラダ・コンスタンティン・ラドゥレス・モトゥル・ヌマルル 4 , ブロック 1 エスセ 2 , アパルタメントゥル 5 8 , セクター 4