

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7661610号
(P7661610)

(45)発行日 令和7年4月14日(2025.4.14)

(24)登録日 令和7年4月4日(2025.4.4)

(51)国際特許分類 F I
G 0 6 F 40/151 (2020.01) G 0 6 F 40/151
G 0 6 F 21/62 (2013.01) G 0 6 F 21/62 3 4 5

請求項の数 15 (全22頁)

(21)出願番号	特願2024-504506(P2024-504506)	(73)特許権者	516005083 ブルー プリズム リミテッド イギリス国 ダブリュエー 2 0 エックス ピー ウォリントン, ファーンヘッド, クラブ レーン, シナモン パーク 2
(86)(22)出願日	令和4年7月29日(2022.7.29)	(74)代理人	110002572 弁理士法人平木国際特許事務所
(65)公表番号	特表2024-530889(P2024-530889 A)	(72)発明者	ジャン, デ イギリス国 シーオー 2 7 エフエイチ コルチェスター, ケンジントン ロード 7 8
(43)公表日	令和6年8月27日(2024.8.27)	(72)発明者	ダッバ, クリシュナ サンディーブ レディ イギリス国 シーエム 6 1 ジービー リ トル キャンフィールド, ワーウィック ロード 1 8
(86)国際出願番号	PCT/EP2022/071384		
(87)国際公開番号	WO2023/012069		
(87)国際公開日	令和5年2月9日(2023.2.9)		
審査請求日	令和6年11月15日(2024.11.15)		
(31)優先権主張番号	21189837.4		
(32)優先日	令和3年8月5日(2021.8.5)		
(33)優先権主張国・地域又は機関	欧州特許庁(EP)		
早期審査対象出願			

最終頁に続く

(54)【発明の名称】 データ難読化

(57)【特許請求の範囲】

【請求項 1】

リモートアクセスアプリケーションを介して受信された機密データがオペレータに対して出力されることを防止するコンピュータ実装方法であって、

リモートアクセスアプリケーションを介してリモートサーバから、前記リモートサーバ上で実行されているソフトウェアアプリケーションのグラフィカルユーザインタフェースを受信するステップ；

前記グラフィカルユーザインタフェースを変更して機密データを削除するステップであって、

プロセッサによって、非構造化画像データを取得するステップ；

前記プロセッサによって、前記非構造化画像データから構造化データを抽出するステップであって、前記構造化データは、機密データであり、定義された機能的フォーマットおよび定義された視覚的フォーマットを有する、ステップ；

前記プロセッサによって、前記構造化データとは異なる人工データを生成するステップであって、前記人工データは、前記構造化データと同じ機能的フォーマットを有する、ステップ；

前記プロセッサによって、前記構造化データが前記人工データに置換された前記非構造化画像データに基づいて人工非構造化画像データを生成するステップであって、前記人工データは、前記構造化データの前記視覚的フォーマットに基づいている、ステップ；

前記プロセッサによって、前記人工非構造化画像データを出力するステップ；

10

20

によって、前記グラフィカルユーザインタフェースを変更して機密データを削除するステップ；

前記変更されたグラフィカルユーザインタフェースを、オペレータが受信するための、コンピュータの1つまたは複数の出力周辺機器に対して出力するステップ；

を有する方法。

【請求項2】

請求項1に記載の方法において、前記人工非構造化画像データは、前記構造化データと同じ視覚的フォーマットを有する方法。

【請求項3】

請求項1または請求項2に記載の方法において、前記非構造化画像データから構造化データを抽出するステップは、

前記非構造化画像データに対して光学式文字認識を実行し、前記非構造化画像データ内のテキストを識別するステップ；

前記テキスト内の構造化データを識別するステップ；

前記構造化データに対応する前記非構造化画像データ内の1つまたは複数のバウンディングボックスを決定するステップ；

前記1つまたは複数のバウンディングボックスを使用して、前記非構造化画像データから1つまたは複数の画像部分を抽出するステップ；

を有する方法。

【請求項4】

請求項3に記載の方法において、前記人工非構造化画像データを生成するステップは、前記1つまたは複数の画像部分における前記構造化データの前記視覚的フォーマットを識別するステップ；

前記1つまたは複数の画像部分における前記構造化データの前記視覚的フォーマットに基づいて、前記1つまたは複数の画像部分に対応する1つまたは複数の人工画像部分を生成するステップ；

前記非構造化画像データを変更して、前記1つまたは複数の画像部分を前記1つまたは複数の人工画像部分によって置換するステップ；

を有する方法。

【請求項5】

請求項1に記載の方法において、前記定義された機能的フォーマットは、エンティティタイプ、およびエンティティタイプフォーマットのうちの1つまたは複数を含む方法。

【請求項6】

請求項1に記載の方法において、前記定義された視覚的フォーマットは、テキスト長、テキストフォント、テキスト色、および背景色のうちの1つまたは複数を含む方法。

【請求項7】

請求項1に記載の方法において、前記人工データを生成するステップは、前記構造化データの機能的フォーマットを識別するステップ；

前記非構造化画像データから抽出された前記構造化データを一覧化するステップ；

前記構造化データの前記機能的フォーマットに基づいて人工データを生成するステップ；

前記人工データを、対応する前記構造化データと共に一覧化するステップ；

を有する方法。

【請求項8】

請求項7に記載の方法において、更に、前記一覧化された構造化データおよび前記一覧化された人工データを編集可能な表として出力することを含む方法。

【請求項9】

10

20

30

40

50

請求項 8 に記載の方法において、更に前記人工非構造化画像データを生成する前に、前記人工データを確認するためのプロンプトをオペレータに送信することを含む方法。

【請求項 10】

請求項 7 記載の方法において、更に、第 2 の非構造化画像データを取得するステップ；前記第 2 の非構造化画像データから第 2 の構造化データを抽出するステップ；前記第 2 の構造化データとは異なる第 2 の人工データを生成するステップを含み、前記第 2 の構造化データが前記構造化データと同一である場合、前記第 2 の人工データは前記人工データと同一である方法。

【請求項 11】

請求項 1 記載の方法において、前記人工非構造化画像データは、前記機密データを使用するソフトウェアの開発プロセス中にディスプレイに出力される方法。

10

【請求項 12】

請求項 11 に記載の方法において、前記ディスプレイは、前記プロセッサとは異なるコンピュータに属する方法。

【請求項 13】

請求項 11 記載の方法において、前記ソフトウェアは、ロボットプロセス自動化 (RPA) プロセスである方法。

【請求項 14】

プロセッサによって実行されると、前記プロセッサに請求項 1 記載の方法を実行させる命令を含む、コンピュータストレージデバイス上に格納されたコンピュータプログラム。

20

【請求項 15】

請求項 1 記載の方法を実行するように構成されたプロセッサを含むコンピューティングシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、機密データを不明瞭化するためのコンピュータ実装方法、コンピュータプログラム及びコンピュータシステムに関する。

【背景技術】

【0002】

データの再編集 (redaction) は、データの機密性が高い内容 (sensitive content)、すなわち認知的内容 (cognitive content) を露出させないように、機密データを視覚的に削除又は置換する処理である。例えば、図 1 A に示す画像には、日付及び / 又は時刻、ドメイン名、電子メールアドレス、氏名、電話番号等、様々な種類の機密データが含まれている。このデータの機密内容は、特定の日付及び / 又は時刻 (Mon 22/02/2021 10:28; 1.5 hrs)、特定のドメイン名 (www.blueprism.com)、特定の電子メールアドレス (Ben.Carter@blueprism.com) 等である。データ再編集とは、このデータを視覚的に除去又は置換する処理である。

30

【0003】

データ再編集を行う理由は、再編集されたデータが共有されるとき、再編集されたデータを受け取った当事者が機密性が高い内容を収集できないようにするためである。機密データの侵害は、機密データの偶発的又は違法な破壊、紛失、改ざん、機密データの不正な開示、又は機密データへの不正なアクセスを引き起こし、重大な人的被害をもたらす可能性があるため、この処理は重要である。更に、受信側当事者による機密データへのアクセスは、様々な法域の規制により違法となる可能性もある。

40

【0004】

データの再編集には幾つかの手法が知られている。これらの手法は、一般に、例えば、不透明な (通常は、黒色の) 矩形、固定文字列 (パスワードの「*****」やクレジットカード番号の「xxxx xxxx xxxx 6789」等)、エンティティタイプに依存する文字列 (人名の「PERSON」等)、又はピクセル化によって、検出された機密データを覆ったり

50

置換したりすることによって、機密データを再編集することを含む。例えば、Google Cloudでは、Cloud Data Loss Preventionは、base 64エンコードされた画像のテキストを検査し、テキスト内の機密データを検出し、照合された機密データが不透明な矩形によって不明瞭化されたbase 64エンコードされた画像を返す。図1Bは、図1Aの画像にこの先行技術を適用したものである。

【発明の概要】

【発明が解決しようとする課題】

【0005】

既存のデータ再編集技術の問題点は、このような技術では、再編集された機密データの視覚的フォーマットのみでなく、機能的フォーマットを理解することも困難になるという10
ことである。これは、機密データがロボットプロセスオートメーションアプリケーション（robotic process automation application）等のソフトウェアアプリケーションで使用され、ソフトウェアアプリケーションの開発及びテストが必要な場合に特に問題となる。この理由は、ソフトウェア開発者は、通常、ソフトウェアアプリケーションが正しくセットアップされ及び機能することを保証するために、機密データの機能的フォーマット、場合によっては、視覚的フォーマットを理解する必要があるからである。

【課題を解決するための手段】

【0006】

本発明は、独立請求項によって定義され、更なる任意の特徴は、従属請求項によって定義される。20

【0007】

本発明の第1の側面では、機密データを不明瞭化するためのコンピュータ実装方法が提供され、この方法は、プロセッサによって、画像データを取得することと、プロセッサによって、画像データから構造化データを抽出することと、構造化データは、機密データであり、定義された機能的フォーマット及び定義された視覚的フォーマットを有することと、プロセッサによって、構造化データとは異なる人工データを生成することと、人工データは、構造化データと同じ機能的フォーマットを有することと、プロセッサによって、構造化データが人工データに置換された画像データに基づいて人工画像データを生成することと、人工データは、構造化データの視覚的フォーマットに基づいて30
いることと、プロセッサによって、人工画像データを出力することとを含む。これにより、データの機密性の高い認知的内容（sensitive cognitive content）を明らかにすることなく、機密データの機能的フォーマットを維持しながら、機密データを不明瞭化できる。これは、（例えば、機能的フォーマットがロボットプロセス自動化（RPA）アプリケーションのような別のソフトウェアアプリケーションとインタラクションするため）機密データの機能的フォーマットを理解する必要があるが、その内容の機密性のために機密データにアクセスできない人間であるオペレータ（human operator：以下、単に「オペレータ」という）にとって有用である。更に、人工データの視覚的フォーマットが機密データの視覚的フォーマットに基づいていることは、人工画像データが画像データの代わりに使用できることを保証するために有用である。

【0008】

幾つかの実施形態では、人工画像データは、構造化データと同じ視覚的フォーマットを有する。これは、オペレータによって開発される必要がある機密データとインタラクションするソフトウェアアプリケーションが、例えば、RPAアプリケーションにおいて、機密データの視覚的フォーマットに何らかの形で依存する場合に有用である。40

【0009】

ある実施形態では、画像データから構造化データを抽出することは、画像データに対して光学式文字認識を実行し、画像データ内のテキストを識別することと、テキスト内の構造化データを識別することと、構造化データに対応する画像データ内の1つ又は複数のバウンディングボックスを決定することと、1つ又は複数のバウンディングボックスを使用して、画像データから1つ又は複数の画像部分を抽出することとを含む。これにより、機50

密データを置換する目的で、機密データを含む画像データの部分を識別し、残りの画像データを無視でき、後続のステップの処理要件を軽減できる。

【0010】

オプションとして、定義された機能的フォーマットは、エンティティタイプ、及びエンティティタイプフォーマットのうちの1つ又は複数を含む。エンティティタイプは、例えば、氏名、日付及び/又は時刻、電子メールアドレス、住所等、機密データが何を示しているかについての情報を提供する。これにより、使用可能なエンティティタイプフォーマットの種類を絞り込むことができる。エンティティタイプフォーマットは、所与のエンティティタイプのフォーマットである。エンティティタイプフォーマットを判定することにより、機密データと同じエンティティタイプフォーマットを使用して人工データを生成できる。例えば、機密データが英国の郵便番号の場合、人工データも英国の郵便番号フォーマットになる。これにより、他のソフトウェアアプリケーションが機密データの機能的フォーマットをどのように扱うか、また、これがそのソフトウェアアプリケーションの実行にどのように影響するかを理解しやすくなる。

10

【0011】

オプションとして、定義された視覚的フォーマットは、テキスト長、テキストフォント、テキスト色、及び背景色のうちの1つ又は複数を含む。テキスト長を使用することにより、生成される人工画像データが、同じ長さ又は類似の長さを使用して、人工データが機密データを囲むバウンディングボックス内に収まるようにできる。テキストフォントを使用することによって、フォントの視覚的フォーマットに関する潜在的な問題を認識できる。例えば、あるソフトウェアアプリケーションが画像上で光学式文字認識を使用しているが、使用されている特定のフォントではI（大文字の「i」）とl（小文字の「L」）の見た目がよく似ている場合、人工画像データ内のテキストフォントを維持することで、そのような問題が識別可能になる。テキスト色と背景色により、人工画像データは、機密データと同じ視覚的フォーマットを持つことができる。

20

【0012】

ある実施形態では、人工データを生成することは、構造化データの機能的フォーマットを識別することと、画像から抽出された構造化データを一覧化することと、構造化データの機能的フォーマットに基づいて人工データを生成することと、人工データを、対応する構造化データと共に一覧化することとを含む。これらの実施形態において、本方法は、オプションとして、一覧化された構造化データ及び一覧化された人工データを編集可能な表として出力することを含んでもよい。これにより、オペレータが人工データを検証し、承認できる。更に重要な点として、オペレータは、編集可能な表を変更して、特定された機密データに修正を加えることができる。具体的には、ユーザは、表に行を追加することによって見逃された機密データを追加し、表の行を削除することによって誤検出された機密データエンティティを削除し、表の対応するセルを編集することによって機密データの誤りを修正し、又は生成された人工データを単に置換することができる。これらの実施形態において、本方法は、更にオプションとして、人工画像データを生成する前に、人工データを確認するためのプロンプトをオペレータに送信することを含んでもよい。これにより、オペレータは、機密データから人工データへのマッピングがどのように実行されるかを透過的に見ることができ、アプローチの堅牢性が向上する。

30

40

【0013】

更なる実施形態において、本方法は、第2の画像データを取得することと、第2の画像データから第2の構造化データを抽出することと、第2の構造化データとは異なる第2の人工データを生成することとを含んでもよく、第2の構造化データが先の構造化データと同一である場合、第2の人工データは先の人工データと同一である。これは、同じ機密データの複数回の出現が、本方法を用いて常に同様に難読化されることを意味する。このように人工データを異なる画像データ間で一貫させることにより、人工データを使用するオペレータに混乱を生じさせることが回避される。更に、データの完全性が維持されるため、機密データを使用するソフトウェアアプリケーションのテスト/デバッグが可能になる。

50

【 0 0 1 4 】

特定の実施形態では、人工画像データを生成することは、1つ又は複数の画像部分における構造化データの視覚的フォーマットを識別することと、1つ又は複数の画像部分における構造化データの視覚的フォーマットに基づいて、1つ又は複数の画像部分に対応する1つ又は複数の人工画像部分を生成することと、画像データを変更して、1つ又は複数の画像部分を1つ又は複数の人工画像部分に置換することとを含んでもよい。このように画像部分を考慮することにより、機密データを不明瞭化するために異なる画像部分における異なる視覚的フォーマットを考慮に入れることができる。

【 0 0 1 5 】

幾つかの実施形態では、人工画像データは、機密データを使用するソフトウェアの開発プロセス中にディスプレイに出力される。データの機能的フォーマットと、データの視覚的フォーマットは、いずれも、データのこれらの側面を何らかの方法で使用する様々なソフトウェアアプリケーションにとって重要であり得る。機能的フォーマットと視覚的フォーマットを維持することにより、ソフトウェア開発者は、そのソフトウェアアプリケーションが正しく設定され、機能していることを確認できる。オプションとして、ディスプレイは、プロセッサとは別のコンピュータに属する。通常、ソフトウェア開発を行う当事者は、機密データを見ることを許可されない。これにより、一方の当事者（すなわち、機密データを見ることを許可された当事者、クライアントコンピュータ）がそのデータの不明瞭化と人工画像データの生成を実行し、他方の当事者（開発者コンピュータ）がその人工画像データを使用することが保証される。

【 0 0 1 6 】

特定の実施形態では、ソフトウェアは、ロボットプロセス自動化（robotic process automation：RPA）プロセスである。RPAは、自動化されたプロセス及びワークフローを使用して、グラフィカルユーザインタフェース（GUI）又はドキュメント等の画像データから構造化データ、典型的には、機密データを抽出することを含む。自動化されたプロセスやワークフローが正しく機能しているかどうかを判断できるようにするためには、人工データ内の機密データの機能的フォーマットを維持することが重要である。RPAは、機密データを抽出するためにソフトウェアアプリケーションのGUIに依存することが多いため、機密データの視覚的フォーマットを維持することも重要である。すなわち、機密データの視覚的フォーマットを維持することで、抽出に伴う潜在的な問題を明らかにできる。

【 0 0 1 7 】

本発明の第2の側面では、リモートアクセスアプリケーションを介して受信された機密データがオペレータに出力されることを防止するコンピュータ実装方法が提供される。この方法は、リモートアクセスアプリケーションを介してサーバから、リモートサーバ上で実行されるソフトウェアアプリケーションのGUIを受信することと、本発明の第1の側面の方法に基づいて、GUIを変更して機密データを除去することと、変更されたGUIを、オペレータが受信するための、コンピュータの1つ又は複数の出力周辺機器に出力することとを含む。これにより、本発明の第1の側面の利点をリモートアクセスアプリケーションのコンテキストで使用でき、これは、特にソフトウェア開発目的に有用である。

【 0 0 1 8 】

本発明の第3の側面では、プロセッサによって実行されると、プロセッサに本発明の第1の側面の方法を実行させる命令を含むコンピュータプログラムが提供される。

【 0 0 1 9 】

本発明の第4の側面では、プロセッサによって実行されると、プロセッサに本発明の第1の側面の方法を実行させる命令を含むコンピュータ可読媒体が提供される。

【 0 0 2 0 】

本発明の第5の側面では、本発明の第1の側面の方法を実行するように構成されたプロセッサが提供される。

【 0 0 2 1 】

10

20

30

40

50

本発明の第 6 の側面では、本発明の第 1 の側面の方法を実行するように構成されたプロセッサを含むコンピューティングシステムが提供される。

【 0 0 2 2 】

以下の図面を参照して、本発明の実施形態を例示的に説明する。

【図面の簡単な説明】

【 0 0 2 3 】

【図 1 A】本発明の方法で使用する画像データの例を示す図である。

【図 1 B】（先行技術）公知のデータ再編集技術を適用した後の図 1 A の画像データを示す図である。

【図 2】本発明の方法を実施するためのシステムの例を示す図である。

10

【図 3】本発明の方法を実施するためのシステムの例を示す図である。

【図 4】本発明の方法に基づいて不明瞭化される機密データを含む画像データの例を示す図である。

【図 5】本発明の方法を示す図である。

【図 6 A】本発明の方法に基づく、図 1 A の画像データ例からの構造化データの抽出を示す図である。

【図 6 B】本発明の方法で使用するための、図 1 A の構造化データに対応するデータベースを示す図である。

【図 6 C】本発明の方法に基づいて図 1 A の構造化データが置換された人工画像データを示す図である。

20

【発明を実施するための形態】

【 0 0 2 4 】

図 2 は、一実施形態に基づいて本発明の方法が実施されるコンピューティングシステム 10 を示している。コンピューティングシステム 10 は、1 人以上のオペレータ 25 (human operator 25) が物理的にアクセス可能な 1 台以上のクライアントコンピュータ 20 を備える。

【 0 0 2 5 】

また、コンピューティングシステム 10 は、1 つ又は複数のサーバ 50 を備える。サーバ 50 は、通常、リモートサーバであり、すなわち、サーバ 50 は、オペレータ 25 が物理的にアクセスできないように、クライアントコンピュータ 20 とは異なる場所に配置されている。幾つかの例では、リモートサーバ 50 は、仮想サーバである。図 2 では、サーバ 50 - 1、50 - 2、50 - 3 は、クラウドコンピューティング環境 60 内の仮想サーバである。クライアントコンピュータ 20 及びサーバ 50 は、少なくとも 1 つの通信ネットワーク 30 を介して互いに通信可能に接続されている。この通信可能な接続により、クライアントコンピュータ 20 とサーバ 50 との間でデータを通信できる。少なくとも 1 つの通信ネットワーク 30 は、通常、インターネット（すなわち、IP、IPv4、IPv6）を含む。インターネットに加えて又はインターネットに代えて、セルラーネットワーク（すなわち、3G、4G LTE、5G）、ローカルエリアネットワーク、クラウドネットワーク、無線ネットワーク、又は任意の他の既知の通信ネットワーク等の他の通信ネットワークが存在してもよい。

30

40

【 0 0 2 6 】

また、コンピューティングシステム 10 には、ソフトウェア開発者であるオペレータ 45 (human operator 45) がアクセス可能な開発者コンピュータ 40 が存在する。開発者コンピュータ 40 は、サーバ 50 に通信可能に接続されており、これにより、オペレータ 45 は、サーバ 50 上で実行されているソフトウェアアプリケーション 50 A をセットアップ、開発、構成、スケジューリング、又は監視できる。これに代えて又はこれに加えて、開発者コンピュータ 40 は、クライアントコンピュータ 20 に通信可能に接続されており、これにより、オペレータ 45 は、クライアントコンピュータ 20 上で実行されているソフトウェアアプリケーション 20 A をセットアップ、開発、構成、スケジューリング、又は監視できる。

50

【 0 0 2 7 】

ここに説明したコンピューティングシステム 10 は、例示的なものに過ぎず、システム構成要素の削除又は追加を含む変更が可能である。

【 0 0 2 8 】

図 3 は、図 2 に示すコンピューティングシステム 10 の選択された側面を示している。具体的には、図 3 は、通信ネットワーク 30 を介してサーバ 50 と通信するクライアントコンピュータ 20 を示している。クライアントコンピュータ 20 は、1つ又は複数のソフトウェアアプリケーション 20 A、プロセッサ 20 B、メモリ 20 C、1つ又は複数の入力周辺機器 20 D、及び1つ又は複数の出力周辺機器 20 E を備える。プロセッサ 20 B は、中央処理装置 (central processing unit : CPU) 及び / 又はグラフィック処理装置 (graphical processing unit : GPU) を含む。メモリ 20 C は、データ記憶装置及び / 又は半導体メモリを含む。データ記憶装置は、ハードディスクドライブ、ソリッドステートドライブ、外部ドライブ、リムーバブル光ディスク、及び / 又はメモリカードの形態をとる。半導体メモリは、データを一時的に記憶するための揮発性メモリ、例えば、ランダムアクセスメモリ (RAM)、及びデータを長期的に記憶するための不揮発性メモリ、例えば、リードオンリメモリ (ROM)、フラッシュメモリの形態をとる。

10

【 0 0 2 9 】

1つ又は複数のソフトウェアアプリケーション 20 A は、メモリ 20 C にコンピュータプログラムとして格納され、クライアントコンピュータ 20 上でプロセッサ 20 B によって実行される。入力周辺機器 20 D 及び出力周辺機器 20 E を介してオペレータ 25 との直接的なインタラクションを実現するこれらのソフトウェアアプリケーションは、オペレーティングシステム (OS) 及びデスクトップアプリケーションを含む。既知のオペレーティングシステムの例には、マイクロソフトウィンドウズ (登録商標)、Mac OS、及び Linux (登録商標) が含まれる。クライアントコンピュータ 20 用の既知のデスクトップアプリケーションの例には、Google Chrome 等のウェブブラウザ、Microsoft Word 等の文書作成アプリケーション、Microsoft によるリモートデスクトッププロトコル (remote desktop protocol : RDP) やリモートフレームバッファ (remote framebuffer : RFB) プロトコル等のリモートアクセスアプリケーションが含まれる。但し、本発明は、ここで述べた特定のアプリケーションと共に使用することに限定されるものではない。

20

30

【 0 0 3 0 】

上述したように、クライアントコンピュータ 20 は、1つ又は複数の入力周辺機器 20 D を備える。入力周辺機器 20 D の目的は、オペレータ 25 がクライアントコンピュータ 20 に指示を送ることを可能にすることである。入力周辺機器 20 D の例には、マウス、キーボード、タッチスクリーン、イメージスキャナ、バーコードリーダー、ゲームコントローラ、マイク、デジタルカメラ、ウェブカメラ等が含まれる。

【 0 0 3 1 】

また、クライアントコンピュータ 20 は、1つ又は複数の出力周辺機器 20 E を備える。出力周辺機器 20 E の目的は、オペレータ 25 がクライアントコンピュータ 20 から情報を受け取れることを可能にすることである。出力周辺機器 20 E の例には、ディスプレイ装置 (例えば、コンピュータモニターやプロジェクタ)、プリンタ、ヘッドフォン、コンピュータスピーカ等が含まれる。入力周辺機器 20 D と同様に、出力周辺機器 20 E は、クライアントコンピュータ 20 と一体化されていてもよく、クライアントコンピュータ 20 の外部に設けられていてもよい。オペレータ 25 は、視覚や聴覚等の感覚を用いてアプリケーション 21 の UI を解釈することにより、出力周辺機器 20 E を用いてクライアントコンピュータ 20 から情報を受け取る。

40

【 0 0 3 2 】

クライアントコンピュータ 20 には、(図 3 には示していない) 他の構成要素も存在する。例えば、コンピュータ 20 は、通信ネットワーク 30 を介した通信を可能にするネットワークアダプタカード、電源、マザーボード、サウンドカード等の1つ又は複数

50

る。

【 0 0 3 3 】

開発者コンピュータ 40 は、図 3 に示すクライアントコンピュータ 20 と同じ構成要素を有する。開発者コンピュータ 40 とクライアントコンピュータ 20 の違いとして、開発者コンピュータ 40 は、サーバ 50 において、サーバ 50 上で実行されている 1 つ又は複数のソフトウェアアプリケーション 50 A をセットアップ、開発、設定、スケジューリング、及び監視できるアクセス権を有しており、クライアントコンピュータ 20 は、このようなアクセス権を有していない。これに加えて又はこれに代えて、開発者コンピュータは、クライアントコンピュータ 20 において、クライアントコンピュータ 20 上で実行されている 1 つ又は複数のソフトウェアアプリケーション 20 A をセットアップ、開発、構成、スケジューリング、及び監視できるアクセス権を有しており、クライアントコンピュータ 20 は、開発者コンピュータ 40 と同等の機能を有していない。したがって、開発者コンピュータ 40 は、オペレータ 45 がこのセットアップ、開発、構成、スケジューリング、及び監視を実行できるように、追加のソフトウェアアプリケーション 20 A を有していてもよい。例えば、マイクロソフト社のリモートデスクトッププロトコル (remote desktop protocol : RDP) やリモートフレームバッファ (remote framebuffer : RFB) プロトコル等のリモートアクセスアプリケーションをこの目的に使用できる。

10

【 0 0 3 4 】

図 3 に示すように、サーバ 50 は、1 つ又は複数のソフトウェアアプリケーション 50 A に加えて、プロセッサ 50 B、メモリ 50 C、及びマシンインタフェース 50 D を備える。1 つ又は複数のアプリケーション 50 A は、メモリ 50 C にコンピュータプログラムとして格納され、プロセッサ 50 B によってリモートサーバ 50 上で実行される。

20

【 0 0 3 5 】

サーバ 50 は、単一のサーバ (例えば、図 2 に示すサーバ 50) 又は複数のサーバ (例えば、図 2 に示すサーバ 50 - 1、50 - 2、50 - 3) の形態をとることができ、あるいは、分散サーバの形態をとることもできる。分散サーバは、構成要素であるコンポーネントに処理とデータを分散して動作する。サーバ 50 は、物理サーバでも仮想サーバでもよい。サーバ 50 が仮想サーバである場合、ソフトウェアアプリケーション 50 A、プロセッサ 50 B、メモリ 50 C 及びマシンインタフェース 50 D は全て、コンピュータシステム 10 のクラウドコンピューティング環境 60 においてホストされる仮想エンティティである。

30

【 0 0 3 6 】

サーバ 50 上の 1 つ又は複数のソフトウェアアプリケーション 50 A は、必ずしも、入力周辺機器 20 D 及び出力周辺機器 20 E を介してオペレータ 25 又はオペレータ 45 と直接インタラクションするとは限らない。これに代えて、1 つ又は複数のソフトウェアアプリケーション 50 A は、通信ネットワーク 30 及びマシンインタフェース 50 D を介してクライアントコンピュータ 20 又は開発者コンピュータ 40 と直接インタラクトするアプリケーションであってもよい。幾つかの例では、サーバ 50 上のソフトウェアアプリケーション 50 A は、オプションとして、開発者コンピュータ 40 上のソフトウェアアプリケーション 20 A を介して、ソフトウェアアプリケーション 50 A のセットアップ、開発、構成、スケジューリング、及び監視を実行する際にオペレータ 45 を支援するための開発者インタフェースを開発者コンピュータ 40 に提供できる。この開発者インタフェースは、クライアントコンピュータ 20 を操作するオペレータ 25 には提供されない。

40

【 0 0 3 7 】

1 つ又は複数のソフトウェアアプリケーション 50 A は、クライアントコンピュータ 20 からのデータ又はクライアントコンピュータ 20 に関連するデータを使用できる。このデータは、本明細書で更に説明する「機密データ」である可能性がある。クライアントコンピュータ 20 からのデータ又はクライアントコンピュータ 20 に関連するデータを使用する例示的なソフトウェアアプリケーション 50 A は、米国特許出願第 14 / 053319 号及び米国特許第 10,469,572 号に記載されているようなロボットプロセス自

50

動化 (robotic process automation : R P A) アプリケーションである。R P A アプリケーションでは、クライアントコンピュータ 2 0 からのデータ又はクライアントコンピュータ 2 0 に関連するデータを使用して、自動化されたプロセスが実行される。

【 0 0 3 8 】

ある特定の実施形態において、コンピュータシステム 1 0 を使用して R P A アプリケーションを実行する場合、図 2 に示すように、クラウドコンピューティング環境 6 0 内の複数の仮想サーバ 5 0 - 1、5 0 - 2、5 0 - 3 に加えて、物理サーバ 5 0 が存在してもよい。本実施形態において、物理サーバ 5 0 は、プロセス定義、ログ、監査、及びユーザ情報を保持する集中リポジトリであるデータベースサーバである。データベースサーバは、複数の仮想サーバのうちの第 1 の仮想サーバ 5 0 - 1 と通信する。第 1 の仮想サーバ 5 0 - 1 は、データベースサーバと、第 2 及び第 3 の仮想サーバ 5 0 - 2、5 0 - 3 との間の接続を制御するアプリケーションサーバである。アプリケーションサーバは、仮想 W i n d o w s サーバとしてプロビジョニングされ、セキュアなクレデンシャル管理 (secure credential management)、データベース接続マーシャリング (database connection marshalling)、データ暗号化、スケジューリングされたプロセス実行等の機能を含んでもよい。第 2 の仮想サーバ 5 0 - 2 は、ロボットプロセス自動化のための自動化されたプロセスの実行を担当する、通常、標準化されたエンドユーザデスクトップの仮想化されたインスタンスをホストする。第 3 の仮想サーバ 5 0 - 3 は、自動化プロセスのセットアップ、開発、設定、スケジューリング、及び監視を実現するエンドユーザデスクトップ構築である。第 3 の仮想サーバ 5 0 - 3 は、専用のソフトウェアアプリケーション 2 0 A を介して、開発者コンピュータ 4 0 にアクセス可能である。

10

20

【 0 0 3 9 】

機密データ

多くの場合、図 2 に示す例示的なコンピュータシステム 1 0 のようなコンピュータシステムは、機密データを扱うことが要求される。例えば、開発者コンピュータ 4 0 は、ソフトウェアアプリケーション 5 0 A のセットアップ、開発、設定、スケジューリング、又は監視を行うために、サーバ 5 0 上で実行されているソフトウェアアプリケーション 5 0 A (例えば、R P A アプリケーション) にアクセスする必要がある場合がある。しかしながら、ソフトウェアアプリケーション 5 0 A は、機密データを使用している可能性があり、すなわち、開発者コンピュータ 4 0 を操作するオペレータ 4 5 によって閲覧されるべきではないクライアントコンピュータ 2 0 からのデータ又はクライアントコンピュータ 2 0 に関連するデータを使用している可能性がある。別の例では、開発者コンピュータ 4 0 は、機密データを使用するクライアントコンピュータ 2 0 上のソフトウェアアプリケーション 2 0 A を設定、開発、構成、スケジューリング、又は監視するために、クライアントコンピュータ 2 0 にアクセスする必要がある場合がある。

30

【 0 0 4 0 】

ここで言う機密データとは、その認知的内容により、高度なセキュリティ配慮を必要とする特別なタイプのデータである。機密データの侵害は、機密データの偶発的又は違法な破壊、紛失、改ざん、機密データの不正な開示、又は機密データへのアクセスを引き起こし、重大な人的被害をもたらす可能性がある。例えば、ある個人の医療記録が永久的に削除されると、その個人の健康に重大かつ長期的な影響が及ぶ可能性がある。このため、欧州連合 (E U) では、一般データ保護規則 (General Data Protection Regulation : G D P R)、英国では、データ保護法 (Data Protection Act 2 0 1 8) 等、様々な法域で機密データの保存と処理が規制されている。

40

【 0 0 4 1 】

機密データは、テキストの形態をとる。機密データには、個人情報、すなわち、特定又は識別可能な自然人に関連する情報が含まれる場合がある。例えば、機密データには、氏名、住所、生年月日、電話番号等が含まれる。他の種類の機密データには、個人の位置データ、オンライン識別情報、個人の身体的、生理的、遺伝的、精神的、経済的、文化的若しくは、社会的アイデンティティに固有の 1 つ又は複数の要素が含まれる。これに加えて

50

又はこれに代えて、機密データには、クレジットカード番号や銀行番号等の金融情報が含まれる場合がある。更に、別の例として、機密データには、医療情報が含まれる場合がある。

【 0 0 4 2 】

本発明の目的のための機密データは、構造化データの形態をとる。本明細書で使用する「構造化データ」という用語は、行（レコード）と列（フィールド）で構造化された、リレーショナルデータベース等の電子ファイル内に格納できるデータを意味する。例えば、日付を取得するには、日付フィールド（日付列）にアクセスする。これに対し、「構造化されていない」自由フォーマットのテキストから意味を導き出すには、テキストを順次スキャンして比較する必要がある。

10

【 0 0 4 3 】

機密データは、定義された視覚的フォーマットのみでなく、定義された機能的フォーマットを有する。本明細書で使用する「機能的フォーマット」という用語は、ソフトウェアアプリケーション 20 A を介してクライアントコンピュータ 20 の動作を制御するような技術システムにおいて技術的機能を有する機密データの部分、特にその部分のフォーマットを指す。ここで使用する「視覚的フォーマット」という用語は、機密データの提示の手法を指す。

【 0 0 4 4 】

機能的フォーマットは、エンティティタイプ及び/又はエンティティタイプフォーマットのうちの1つ又は複数を含むことができる。エンティティタイプは、例えば、氏名、日付及び/又は時刻、電子メールアドレス、住所等、機密データが何を示しているかについての情報を提供する。電子メールアドレスの機能的フォーマットは、電話番号の機能的フォーマットとは異なるため、使用可能なエンティティタイプフォーマットの種類を絞り込むことができる。エンティティタイプフォーマットは、所与のエンティティタイプのフォーマットである。日付及び/又は時刻のエンティティタイプは、様々なエンティティタイプフォーマットを有することができる。例えば、日付は、「DD/MM/YYYY」（欧州標準フォーマット）、「MM/DD/YYYY」（米国標準フォーマット）、「YYYY/MM/DD」（日本標準フォーマット）で表すことができる。「D」は、日、「M」は、月、「Y」は、年を表す。時間は「hh:mm」、「hh:mm:ss」で表すことができ、「h」は、24時間制の時、「m」は、分、「s」は、秒である。12時間制（AM、PM）、異なるタイムゾーン（GMT、EST等）等、その他の時間フォーマットもある。別の例において、英国の郵便番号の場合、エンティティタイプのフォーマットは、「A9AA 9AA」、「A9A 9AA」、「A9 9AA」、「A99 9AA」、「AA9 9AA」、又は「AA99 9AA」（ここで、「A」は文字、「9」は、数字を表す）といった設定された数のフォーマットのいずれかでなければならない。

20

30

【 0 0 4 5 】

視覚的フォーマットは、テキスト長、テキストフォント、テキスト色、及び/又は背景色の1つ又は複数を含んでいてもよい。テキスト長は、テキストの水平方向の長さを表す。テキスト長は、テキストの文字数に基づいて測定してもよく（例えば、「Ben Carter」は、スペースを含めて10文字）、ピクセル数を使用して測定してもよい。テキストフォントには、テキストの書体（例えば、Arial、Times New Roman、Courier New）、テキストのサイズ（例えば、12pt）、及び特殊なスタイル特性（太字、斜体、下線、取り消し線、下付き文字、上付き文字等）が含まれる。テキスト色は、テキストの主要な色であり、通常はRGBスケール又はHSLスケールで測定される。背景色は、テキストを囲む背景の主な色であり、これも通常はRGBスケール又はHSLスケールで測定される。

40

【 0 0 4 6 】

図4は、機密データ410を含むグラフィカルユーザインタフェース（GUI）400の例である。この例では、機密データ410は、日付及び/又は時刻410A、ドメイン名410B、電子メールアドレス410C、氏名410D、及び電話番号410Eといっ

50

たエンティティタイプを含む。日付及び/又は時刻 4 1 0 Aとして示されるフォーマットには、「ddd DD/MM/YYYY hh:mm」及び「小数時間」の2つのエンティティタイプフォーマットがある。視覚的フォーマットとして、機密データは、同じフォント(Calibri)を有し、グレーのテキスト色と白の背景色を有する。この例については、本発明の方法を説明する文脈で、図6を参照して更に詳しく後述する。このGUIは、RPAアプリケーションに見られる典型的なGUIである。

【0047】

図4に示すタイプのような機密データを再編集するための様々な手法が知られている。しかしながら、このような手法は、クライアントコンピュータ20からの又はクライアントコンピュータ20に関連する機密データがソフトウェアアプリケーション50Aによって使用され、ソフトウェアアプリケーション50Aが開発者コンピュータ40上でオペレータ45によって更に開発及びテストされる必要がある場合のシナリオでの使用には適していない。この理由は、既知の技術では、機密データの機能的フォーマットのみでなく視覚的フォーマットも置換又は除去され、データの機能的フォーマット及び視覚的フォーマットが使用可能にならないため、オペレータ45がソフトウェア開発及びテストタスクを実行することが困難になるからである。

【0048】

方法の概説

図5は、既知のデータ再編集技術と同様の欠点を有することなく、機密データを不明瞭化するために本発明が採用する方法500を示している。図5の方法500は、クライアントコンピュータ20のプロセッサ20Bによって実行してもよく、サーバ50のプロセッサ50Bによって実行してもよい。開発者コンピュータ40は、機密データを含む画像データにアクセスすべきではないため、図5の方法は、開発者コンピュータ40のプロセッサ20Bによって実行されるべきではない。

【0049】

図5に示すように、方法500は、プロセッサによって実行される以下のステップを含む。

【0050】

- ・プロセッサによって、画像データを取得する(ステップ510)
- ・プロセッサによって、機密データであり、定義された機能的フォーマット及び定義された視覚的フォーマットを有する構造化データを画像データから抽出する(ステップ520)
- ・プロセッサによって、構造化データと同じ機能的フォーマットを有する、構造化データとは異なる人工データを生成する(ステップ530)；
- ・プロセッサによって、構造化データが、構造化データと同じ視覚的フォーマットに基づいている人工データに置換された画像データに基づいて人工画像データを生成する(ステップ540)
- ・プロセッサによって、人工画像データを出力する(ステップ550)。

【0051】

方法500は、データの機密性の高い認知的内容を明らかにすることなく、機密データの機能的フォーマット、並びに視覚的フォーマットを維持しながら、機密データを不明瞭化することを可能にする。これは、開発者コンピュータ40のオペレータ45が、クライアントコンピュータ20に関連する又はクライアントコンピュータ20からの機密データを使用するサーバ50におけるソフトウェアアプリケーション50A上のソフトウェア開発及びテストタスク、又はクライアントコンピュータ20に関連する又はクライアントコンピュータ20からの機密データを使用するクライアントコンピュータ20におけるソフトウェアアプリケーション20Aのためのソフトウェア開発及びテストタスクを実行するために特に有用である。

【0052】

ステップ510からステップ550については、以下のセクションで更に詳細に説明す

る。

【 0 0 5 3 】

画像データ取得

図5のステップ510において、プロセッサ20B又はプロセッサ50Bは、画像データを取得する。画像データは、クライアントコンピュータ20に関連付けられ又はクライアントコンピュータ20から取得される。「関連付けられる」とは、画像データがクライアントコンピュータ20のクライアントに属するか又はクライアントコンピュータ20のクライアントに関連することを意味する。「から」とは、クライアントコンピュータ20から間接的に受信される画像データ及び/又はクライアントコンピュータ20から直接的に受信される画像データの両方を意味する。

10

【 0 0 5 4 】

画像データには、クライアントコンピュータ20から見て機密データであるデータが含まれている可能性がある。これは、クライアントコンピュータ20を操作するオペレータ25は、機密データを閲覧し、インタラクトする権限を有するが、開発者コンピュータ40を操作するオペレータ45は、その権限がないことを意味する。

【 0 0 5 5 】

画像データは、グラフィカルユーザインタフェース(GUI)の形態をとることができる。例えば、画像データは、(サーバ50がオペレータによってアクセス可能なGUIを有する実施形態では)クライアントコンピュータ20又はサーバ50からキャプチャされたGUI画像であってもよい。GUIは、クライアントコンピュータ20上で実行されているソフトウェアアプリケーション20A又はサーバ50上で実行されているソフトウェアアプリケーション50Aのうちの1つ又は複数を示すことができる。GUIは、(オペレータに出力されるような)デスクトップインタフェース全体を含んでもよく、ソフトウェアアプリケーション20A又は50Aの特定の1つに関連するデスクトップインタフェースの部分のみを含んでもよい。画像データがGUIの形態をとる場合、プロセッサ20B又はプロセッサ50Bは、GUIをキャプチャすることによって画像データを取得する。これに代えて、GUIは、プロセッサ20B又はプロセッサ50Bによって先にキャプチャされ、それぞれメモリ20C又はメモリ50Cに格納されていてもよい。このような例では、GUIは、メモリ20C又はメモリ50Cから読み出すことができる。幾つかの例では、GUIは、リモートアクセスアプリケーションによってキャプチャされてもよい。

20

30

【 0 0 5 6 】

これに代えて、画像データは、PDF又は画像ファイル等の文書の形態をとってもよい。この場合、プロセッサ20B又はプロセッサ50Bは、画像データが格納されているメモリ20C又はメモリ50Cからそれぞれ画像データを取得する。

【 0 0 5 7 】

図6Aは、画像データの一例を示す図である。具体的には、図6Aは、Microsoft Outlook等の電子メールアプリケーションのGUI600を示している。このGUIには、氏名、日付、電話番号、電子メールアドレス、ドメイン名等の機密データが含まれている。

40

【 0 0 5 8 】

構造化データ抽出

図5のステップ520では、プロセッサ20B又はプロセッサ50B(ステップ510で画像データを取得した一方)が画像データから構造化データを抽出する。上述したように、構造化データは、定義された機能的フォーマット及び定義された視覚的フォーマットを有する機密データである。したがって、このステップの目的は、機密データである画像データ内の構造化データを特定して抽出することである。

【 0 0 5 9 】

画像データは生来的に構造化されていないが、構造化されていない画像データから構造化データを識別する方法は既知である。例えば、Googleは、コンピュータビジョン

50

(光学式文字認識(OCR)を含む)及び自然言語処理(NLP)を使用して、文書に対して事前に訓練されたモデルを作成するDocument AIを使用している。

【0060】

ステップ520を実行するための1つの方法は、最初に画像データに対して光学式文字認識を実行して、画像データ内のテキストを識別し、テキスト内の構造化データを識別することである。これは、Document AI等の既知の方法を使用して実行できる。構造化データがテキスト内で識別された後、ステップ520を実行することは、構造化データに対応する画像データ内の1つ又は複数のバウンディングボックスを決定することと、1つ又は複数のバウンディングボックスを使用して画像データから1つ又は複数の画像部分を抽出することとを含んでもよい。機密データであると識別された構造化データのみについて、画像データから対応する画像部分を抽出することが望ましい。

10

【0061】

上述した方法を用いて抽出された、機密データに対応する構造化データをそれぞれ含む画像データからの1つ又は複数の画像部分に加えて、ステップ520は、第2の出力を有する。すなわち、1つ又は複数の画像部分内の機密データに対応する基底にある構造化データも出力される。このデータは、ステップ530の準備として一覧化(tabulated)される(表にされる)。

【0062】

図6Aの画像データ例であるGUI600は、構造化データを含む複数のバウンディングボックス610を示している。複数のバウンディングボックス610は、機密データを含む複数の画像部分を形成するように抽出される。更に、画像部分内の構造化データが抽出される。GUI600の場合、構造化データには、「Ben Benjamin Carter」、「Ben Carter」、「De Zhang」、「Eric Tyree」、「John Reid」、「Krishna Dubba」、及び「+44 785 407 9884」が含まれる。この構造化データは、図6Bの表620の列615に示すように、一覧化できる。

20

【0063】

人工データ生成

図5のステップ530において、プロセッサ20B又はプロセッサ50B(ステップ510において画像データを取得した一方)は、構造化データとは異なるが、構造化データと同じ機能的フォーマットを有する人工データを生成する。このステップの目的は、機密性の高い内容を含むことなく機能的フォーマットを維持した、構造化データの代替データを提供することである。

30

【0064】

構造化データとは異なるが、構造化データと同じ機能的フォーマットを有する人工データを生成できるようにするには、データの機能的フォーマットを識別できるようにする必要がある。これは、機械学習モデルやヒューリスティックルールを利用することで達成される。機械学習モデルは、ニューラルネットワーク等の機械学習技術を使用して事前に訓練される。具体的には、テキスト中のエンティティタイプを認識するためのニューラルネットワークモデルは、RoBERTa等のトランスフォーマーアーキテクチャに基づいてもよい。例えば、Hugging Face社(<https://huggingface.co/roberta-base>参照)が提供する事前訓練済みRoBERTaモデルの基本バージョン(「roberta-base」)は、768の隠れベクトルサイズを持つ12個のエンコーダ層、12個のアテンションヘッド、125Mのパラメータで構成されている。このモデルは、強化されたBERT(Bidirectional Encoder Representations from Transformers)アルゴリズムを用いて、様々なサイズとドメインの5つの一般公開コーパスから得た160GB以上の英語テキストデータで学習されている。エンティティ認識の目的のために、事前に学習されたroberta-baseモデルは、よく見られる18のエンティティタイプを含むOntoNotes5データセット(<https://deepai.org/dataset/ontonotes-v5-english>等参照)等のラベル付きNER(Named Entity Recognition)データセットによって微調整できる。ヒューリスティックルールには、(重み付けされた)正規表

40

50

現のテキストパターンや、周囲の単語等の文脈上の手がかりを含めることができる。例えば、正規表現「\b([0-9]{10})\b」を使用して、米国の電話番号又は米国の銀行口座番号を表すと思われる10桁のシーケンスをテキストから抽出できる。次に、ローカルコンテキスト（例えば、10桁のシーケンスの前後5語）に「mobile」や「call」等の単語があれば、その10桁のシーケンスは、米国の電話番号である可能性が高い。逆に、ローカルコンテキストに「savings」や「debit」といった単語があれば、10桁の数字列が米国の銀行口座番号である可能性が高い。更に、正規表現を使用してデータの機能的フォーマットを判定することもできる。例えば、正規表現「(\+[0-9]{1,3})?([0-9]{10})\b」を使用すると、10桁の電話番号の前に対応する国の通話コード（プラス記号の後に1~3桁の数字が続き、その後スペース文字が続くフォーマット）を有するかどうかを判断できる。

10

【0065】

エンティティタイプ及びエンティティタイプフォーマットの表示は、メモリ20B又はメモリ50Bに出力又は格納できる。例えば、図6Aでは、バウンディングボックスの1つに「Mon22/02/2021 10:28」という構造化データが含まれている。上述した方法を用いると、この構造化データエンティティタイプは「日付及び/又は時刻」であるとみなされ、エンティティタイプのフォーマットは「ddd DD/MM/YYYY hh:mm」であるとみなされる。

【0066】

次に、人工データを生成するために、オペレータ25が編集可能なテーブルを使用することが好ましい。具体的には、構造化データの機能的フォーマットを特定した後、画像から抽出された構造化データを一覧化する。この一覧化は、識別された機能的フォーマットの補助によって行うことができる。更に、機能的フォーマット自体を一覧化することもできる。図6Bは一例を示しており、表620は、識別された機能的フォーマットエンティティタイプを使用する第1の列625と、対応する構造化データを示す第2の列615とを有する。この例では、エンティティタイプは、アルファベット順にソートされ、これにより、エンティティタイプがグループ化されている。

20

【0067】

その後、テーブルの各行に対して、構造化データと同じ機能的フォーマットを有するが、異なる認知的内容を有し、したがって機密データではない人工データが生成される。例えば、日付及び/又は時刻のエンティティタイプである構造化データ「Mon 22/02/2021 10:28」に対して、これを同じ機能的フォーマットの人工的な日付及び/又は時刻、例えば「Fri 24/06/1987 19:03」に変更できる。その後、人工データは、構造化データ及びオプションとして機能的フォーマットとともに一覧化される。図6Bは、生成された人工データを含む表620の第3列660を示している。

30

【0068】

構造化データと同じ機能的フォーマットを有する人工データの生成は、エンティティタイプとエンティティタイプフォーマットを条件とするエンティティの確率分布からのランダムサンプリングによって実行される。例えば、人工的な人名を生成するには、姓(surname)のリストから姓を、名(first-name)のリストから名をランダムにサンプリングする。そして、姓と名は、構造化データの同じフォーマット（例えば、「SURNAME, First-name」）にまとめられる。別の例として、人工的な日付を生成するために、元の日付から5年以内の距離にあるカレンダーの日付がランダムにサンプリングされ、その後、構造化データの同じフォーマット（例えば、「ddd DD/MM/YYYY」）を使用して、日付のテキスト表現が生成される。オプションとして、生成された人工データが、対応するソフトウェアアプリケーションによって要求される有効範囲内にあることを保証するために、ポストフィルタリングを実施してもよい。

40

【0069】

一覧化された構造化データ及び一覧化された人工データは、オペレータ25が編集できるように、編集可能な表としてクライアントコンピュータ20上に出力してもよい。これ

50

により、オペレータは、構造化データと人工データとのマッピングを適切と思われるように追加、削除、又は編集できる。例えば、オペレータ 25 は、画像データ内の構造化データ検出に修正を加えるように編集可能テーブルを編集できる。具体的には、ユーザは、表に行を追加することによって、見逃された機密データエンティティを追加してもよく、表の行を削除することによって、誤検出された機密データを削除してもよく、表の対応するセルを編集することによって、機密データ検出の誤りを修正してもよく（例えば、図 6 B において、誤って識別された「Ben Benjamin Carter」を「Benjamin Carter」に変更してもよく）、又は自動的に生成された偽のデータエンティティを、ユーザが好むものに単に置換してもよい（例えば、図 6 B において、Krishna Dubba の偽の氏名を「Michele Oneal」を「Dr Strangelove」に変更してもよい）。

10

【0070】

ステップ 540 において人工画像データを生成する前にオペレータ 25 に人工データを承認させることは有用である。したがって、プロセッサ 20 B は、人工データを確認するためのプロンプトをオペレータ 25 に送信してもよい。

【0071】

編集可能テーブルは、異なる（例えば、第 2、第 3 等の）画像データからの構造化データ及び人工データを含むことができる。このような例では、後続の構造化データ（例えば、第 2 の構造化データ）が先の構造化データと同一である場合、後続の画像データ（例えば、第 2 の人工データ）に基づいて生成される人工データは先の人工データと同一であることが好ましい。例えば、図 6 A の GUI 600 を参照すると、第 2 の画像データが、Ben Carter からの電子メールを含む更なる GUI に関連する場合、構造化データ「Ben Carter」は、元の画像データ及び後続の第 2 の画像データの両方に現れることになる。そして、構造化データ「Ben Carter」が元の画像データに由来するか第 2 の画像データに由来するかにかかわらず、生成される人工データは、Anne Wells となる（図 6 B 参照）。編集可能テーブルは、第 2 の画像データの「Ben Carter」のインスタンスのために新たな行を生成する必要はなく、元の画像データから既に生成されている行を利用できる。このように、編集可能テーブルは、1 つ又は複数の画像データにわたる構造化データと人工データとのマッピングのグローバルテーブルであるとみなすことができる。

20

【0072】

人工画像データの生成

30

図 5 のステップ 540 において、プロセッサ 20 B 又はプロセッサ 50 B は、構造化データが人工データに置換された画像データに基づいて人工画像データを生成し、人工データは、構造化データと同様の視覚的フォーマットに基づいている。

【0073】

ステップ 540 を実行するために、プロセッサ 20 B 又はプロセッサ 50 B は、まず、1 つ又は複数の画像部分における構造化データの視覚的フォーマットを識別し、次に、1 つ又は複数の画像部分における構造化データの視覚的フォーマットに基づいて、1 つ又は複数の画像部分に対応する 1 つ又は複数の人工画像部分を生成し、次に、画像データを変更して 1 つ又は複数の画像部分を 1 つ又は複数の人工画像部分に置換する。

【0074】

40

視覚的フォーマットの識別は、Python Image Library (Pillow) に実装されているようなデジタル画像処理技術や、OpenCV ライブラリによって提供されているようなコンピュータビジョン技術を利用することによって実行される。テキスト長、テキスト色、背景色を検出するための既知の技術が存在する。テキストフォントの認識は、フォントの書体、フォントサイズ、及びフォントスタイルの可能な構成に対してグリッド検索又はベイズ最適化を実行し、どの構成が元の画像部分に最も類似した構造化データの画像を生成するかを見つけることによって達成できる。

【0075】

1 つ又は複数の人工画像部分の生成は、Python Image Library (Pillow) に実装されているようなデジタル画像処理技術を利用して実行される。例え

50

ば、Python Image Library (Pillow) は、新しい画像を作成したり、既存の画像に注釈を付けたり、レタッチしたりするために使用できる ImageDraw モジュールを提供する。具体的には、ImageDraw.rectangle 関数を使用すると、指定された位置に、指定されたサイズで、指定された塗りつぶし色を背景とする矩形ボックスを描画でき、ImageDraw.text 関数を使用すると、指定された位置に、指定されたフォントと色で、指定されたテキスト（すなわち、元の機能的フォーマットで生成された人工データ）の一部を描画できる。ここで、テキストフォント、テキスト色、及び背景色は、全て、上述の方法を介して識別された視覚的フォーマットによって与えられる。

【0076】

幾つかの例では、人工画像データは、画像データ中の構造化データと同じ視覚的フォーマットを有するように生成される。これにより、構造化データの視覚的フォーマットに依存するソフトウェアアプリケーション 20A 上でソフトウェア開発を行うオペレータ 45 が視覚的フォーマットが何であるかを高い精度で理解できるようになる。あるいは、人工画像データは、人工データに関係する部分と画像データに関係する部分とが区別されるように、画像データの構造化データとは 1 つ又は複数の相違点を有する視覚的フォーマットを有するように生成できる。例えば、図 6C の人工 GUI 650 の人工画像データにおいて、人工画像部分は、これらの部分が人工的であることを示すために、強調表示された背景を有する。この場合でも、人工画像データは、構造化データと同様の視覚的フォーマットに基づくものとみなされる。

【0077】

人工画像データの出力

図 5 のステップ 550 において、プロセッサ 20B 又はプロセッサ 50B（画像データを受信した一方）は、人工画像データを出力する。

【0078】

人工画像データは、メモリ 20C 又はメモリ 50C に格納されるファイルに出力してもよい。このファイルは、後日、開発者コンピュータ 40 において、オペレータ 45 がソフトウェア開発の目的でアクセスできる。ファイルは、人工画像データのみを含んでもよい。このような場合、ファイルは、元の画像データと同じファイルフォーマット（例えば、pdf、jpeg）であってもよい。あるいは、ファイルは、人工画像データと、ソフトウェア開発をガイドしサポートするための他の関連情報とを含む包括的なドキュメント（comprehensive document）であってもよい。このようなファイルの例としては、ロボットプロセス自動化（robotic process automation：RPA）内で開発されるビジネスプロセスのフローをキャプチャするプロセス設計文書（process design document：PDD）が挙げられる。

【0079】

これに代えて又はこれに加えて、人工画像データは、ディスプレイに出力してもよい。例えば、オペレータ 25 がチェックするために、人工画像データをクライアントコンピュータ 20 に出力してもよい。他の例では、人工画像データは、プロセッサ 20B 又はプロセッサ 50B とは異なるコンピュータに属するディスプレイに出力してもよい。具体的には、開発者コンピュータ 40 のディスプレイに人工画像データを出力してもよく、これにより、機密データを使用するソフトウェアのソフトウェア開発プロセス中に、オペレータ 45 が画像データ内の機密データを閲覧できないようにできる。

【0080】

幾つかの例では、人工画像データを使用して、機密データを使用するロボットプロセス自動化（RPA）アプリケーションを開発してもよい。RPA では、自動化されたプロセスやワークフローを使用して、GUI やドキュメント等の画像データから、通常は、機密データである構造化データを抽出する。自動化されたプロセスやワークフローが正しく機能しているかどうかを判断できるようにするためには、人工データ内の機密データの機能的フォーマットを維持することが重要である。RPA は、機密データを抽出するためにソ

10

20

30

40

50

ソフトウェアアプリケーションの GUI に依存することが多いため、機密データの視覚的フォーマットを維持することも重要である。すなわち、機密データの視覚的フォーマットを維持することで、抽出に関する潜在的な問題を明らかにできる。

【0081】

幾つかの例では、人工画像データは、Microsoft によるリモートデスクトッププロトコル (remote desktop protocol: RDP) やリモートフレームバッファ (remote framebuffer: RFB) プロトコル等のリモートアクセスアプリケーションで使用されることがある。例えば、米国特許出願 US 17/144,640 号 (Method of Remote Access) で使用されている GUI 画像変更技術の代わりに、図 5 の方法 500 と出力された人工画像データを使用してもよい。これにより、リモートアクセスプロトコルを介して受信された機密データがオペレータ (例えば、オペレータ 45) に出力されることを防止するコンピュータ実装方法が提供される。このような例では、この方法は、リモートアクセスアプリケーションを介してサーバ 50 から、リモートサーバ上で実行されているソフトウェアアプリケーション 50A の GUI を受信することと、本発明の方法に基づいて、機密データを不明瞭化するために GUI を変更することと、変更された GUI を、オペレータ 45 が受信するための、開発者コンピュータ 40 の 1 つ又は複数の出力周辺機器 40E に出力することとを含む。

【0082】

総括

ソフトウェアで実施する場合、本発明は、コンピュータプログラムの形態をとることができる。コンピュータプログラムは、プロセッサによって使用され又はプロセッサに関連して使用されるコンピュータ実行可能コードを有するコンピュータ可読媒体として具現化できる。コンピュータ可読媒体は、プロセッサによって使用され又はプロセッサに関連して使用されるプログラムを格納、保存、通信、伝播、又は輸送できる任意の有形デバイスである。更に、コンピュータ可読媒体は、電子、磁気、光学、電磁、赤外線、半導体デバイス、又は伝搬媒体であってもよい。コンピュータ読取可能媒体の例としては、半導体メモリ、ランダムアクセスメモリ (random access memory: RAM)、読取専用メモリ (read-only memory: ROM)、フラッシュメモリ、ハードディスクドライブ、ソリッドステートドライブ、光ディスク、メモリカード等がある。現在の光ディスクの例としては、CD、DVD、ブルーレイ等がある。現在のメモリカードの例としては、USB フラッシュドライブ、SD カード、microSD カード、MMC カード、xD カード、メモリスティック等がある。

【0083】

ハードウェアで実施する場合、本発明は、本明細書で説明する特定のハードウェアに限定されない。本発明は、図 2 及び図 3 に関して説明したものと異なるハードウェアで実装しても、上述のように機能させることができることは、当業者にとって明らかである。

【0084】

図示したフローチャートは、本発明の方法の可能な実施側面のアーキテクチャ、機能、及び動作を示している。幾つかの代替的な具体例では、図に示すステップは、図に示す順序とは異なる順序で実行してもよい。例えば、連続して示されている 2 つのステップは、実際には、実質的に同時に実行してもよく、ステップを表すブロックは、関係する機能によっては、逆の順序で実行してよいこともある。

【0085】

上の説明は、例示のみを目的とし、当業者によって様々に変更できることは明らかである。上では、ある程度の特殊性をもって又は 1 つ又は複数の個々の実施形態を参照して様々な実施形態を説明したが、当業者は、本発明の範囲から逸脱することなく、ここに開示した実施形態に多数の変更を加えることができる。

10

20

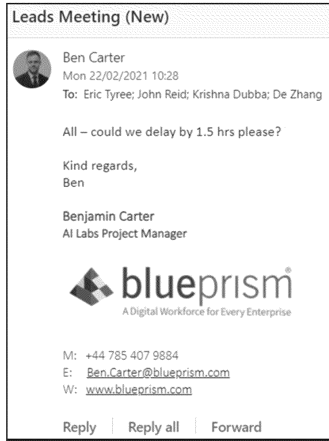
30

40

50

【 面 】

【 1 A 】



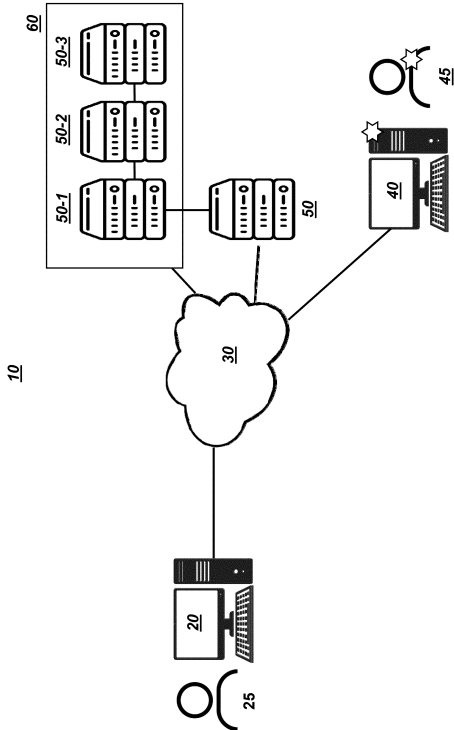
【 1 B 】



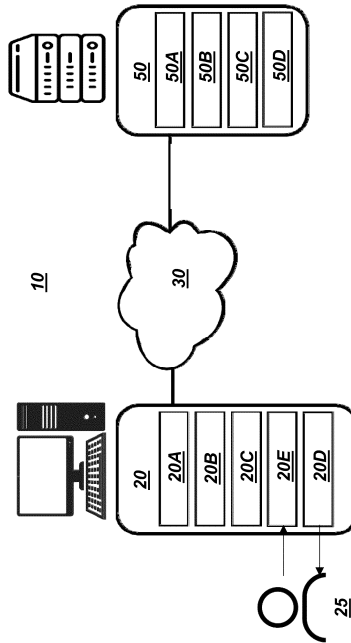
10

20

【 2 】



【 3 】

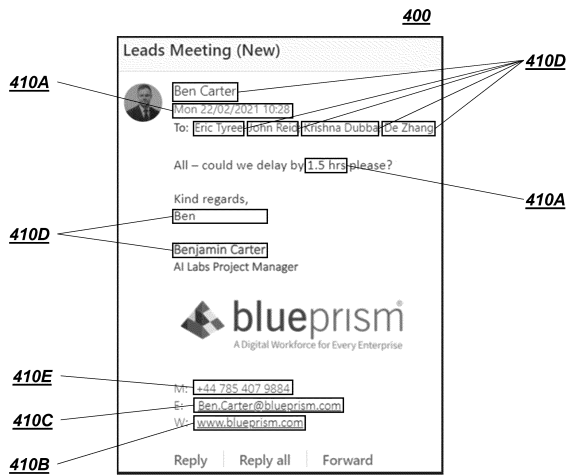


30

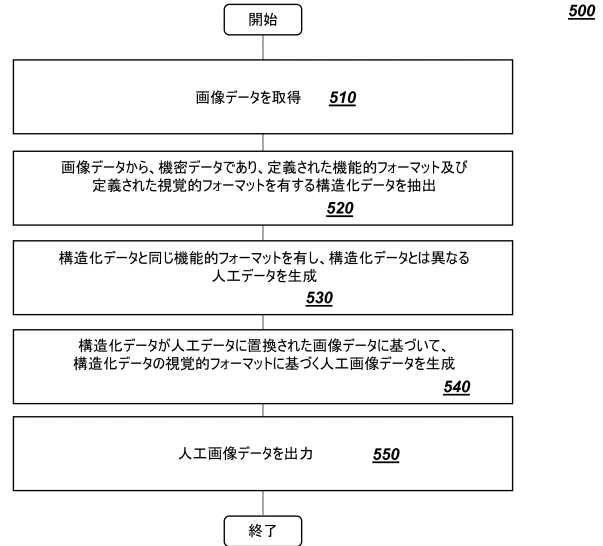
40

50

【 図 4 】

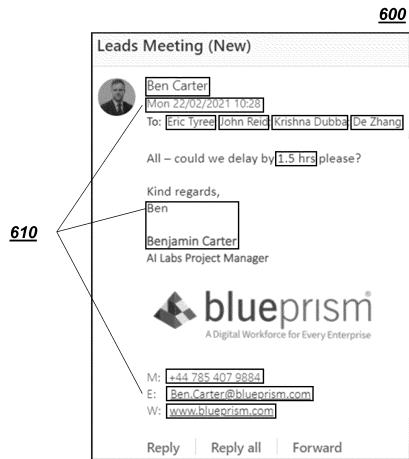


【 図 5 】



10

【 図 6 A 】



【 図 6 B 】

620

エンティティタイプ	625	構造化データ	615	人工データ	660
DATE_TIME	1.5 hrs	1.22 hrs			
DATE_TIME	Mon 22/02/2021 10:28	Mon 22/02/2021 10:28			Fri 24/06/1987 19:03
DOMAIN_NAME	www.blueprism.com	www.blueprism.com			Baker-stewart.com
EMAIL_ADDRESS	Ben.Carter@blueprism.com	Ben.Carter@blueprism.com			williamstonya@powell.net
PERSON	Ben Benjamin Carter	Ben Benjamin Carter			Christopher Farrell
PERSON	Ben Carter	Ben Carter			Anne Wells
PERSON	De Zhang	De Zhang			Ian White
PERSON	Eric Tyree	Eric Tyree			Lisa Mason
PERSON	John Reid	John Reid			Dale Cook
PERSON	Krishna Dubba	Krishna Dubba			Michele Oneal
PHONE_NUMBER	+44 785 407 9884	+44 785 407 9884			978-821-5337x022

30


40

50

【 ☒ 6 C 】


650


Leads Meeting (New)

 **Anne Wells**
Fri 24/06/1987 19:03

To: [Lisa Mason](#) [Dale Cook](#) [Michele Oneal](#) [Ian White](#)

All – could we delay by **1.22 hrs** please?

Kind regards,

Christopher Farrell
AI Labs Project Manager

 **blueprism**
A Digital Workforce for Every Enterprise

M: [978-821-5337x022](tel:978-821-5337x022)
E: williamstonva@powell.net
W: baker-stewart.com

Reply | Reply all | Forward

660

10

20

30

40

50

フロントページの続き

審査官 石坂 知樹

- (56)参考文献 国際公開第2020/082187(WO, A1)
米国特許出願公開第2018/0285591(US, A1)
英国特許出願公開第02531713(GB, A)
- (58)調査した分野 (Int.Cl., DB名)
G06F 40/151
G06F 21/62